



## King's Research Portal

DOI:

[10.1016/j.jeem.2019.05.001](https://doi.org/10.1016/j.jeem.2019.05.001)

*Document Version*

Peer reviewed version

[Link to publication record in King's Research Portal](#)

*Citation for published version (APA):*

Baldwin, E., Cai, Y., & Kuralbayeva, K. (2020). To Build or not to Build? Capital stocks and climate policy. *JOURNAL OF ENVIRONMENTAL ECONOMICS AND MANAGEMENT*, 100, [102235].  
<https://doi.org/10.1016/j.jeem.2019.05.001>

### **Citing this paper**

Please note that where the full-text provided on King's Research Portal is the Author Accepted Manuscript or Post-Print version this may differ from the final Published version. If citing, it is advised that you check and use the publisher's definitive version for pagination, volume/issue, and date of publication details. And where the final published version is provided on the Research Portal, if citing you are again advised to check the publisher's website for any subsequent corrections.

### **General rights**

Copyright and moral rights for the publications made accessible in the Research Portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognize and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the Research Portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the Research Portal

### **Take down policy**

If you believe that this document breaches copyright please contact [librarypure@kcl.ac.uk](mailto:librarypure@kcl.ac.uk) providing details, and we will remove access to the work immediately and investigate your claim.

# To Build or Not to Build?

## Capital Stocks and Climate Policy\*

Elizabeth Baldwin<sup>†</sup>   Yongyang Cai<sup>‡</sup>   Karlygash Kuralbayeva<sup>§</sup>

May 3, 2019

### Abstract

We investigate (i) the impact of emission reduction policy on investment in polluting infrastructure, such as coal-fired power stations and (ii) optimal subsidies for “clean” alternatives with “learning” spillovers. We build a general theoretical model, and embed it in a fully calibrated integrated assessment model. Because emission reduction policy reduces investments in polluting assets, short-term emission reductions are enhanced—our “irreversibility effect”. Thus, “stranded assets” in this fuel-*using* sector have distinctive properties. We also provide a simple formula for how the optimal subsidy to deployment of a “clean” sector depends on its rate of “learning-by-doing” and on its socially-optimal growth. So, if the sector should grow faster for other reasons, its optimal subsidy is increased, showing that its optimal growth rate is faster still—our “acceleration effect”. Our calibrations show that, to limit global climate change to 2°C warming, investments in coal-fired power stations must end very soon. Considering second-best settings, we show that carbon taxes achieve stringent policy targets more efficiently, but subsidies to the “clean” sector deliver higher welfare, and are more efficient, when policy targets are more mild.

**Keywords:** Infrastructure; Clean and Dirty Energy Inputs; Renewable Energy; Stranded Assets; Carbon Budget; Climate Policies; Green Paradox

**JEL codes:** O44, Q54, Q58

---

\*We thank Lassi Ahlvik, Alex Bowen, Maria Carvalho, Simon Dietz, Carolyn Fischer, Roger Fouquet, Reyer Gerlagh, Kenneth Gillingham, Niko Jaakkola, Paul Klemperer, Per Krusell, Linus Mattauch, Armon Rezai, Daniel Spiro, Rick van der Ploeg, Till Requate, Andreas Tryphonides, Frank Venmans, and three anonymous referees and participants at: the Sustainable Development workshop in Rimini (2017); GTAP 2017 (Purdue); the pre-EAERE 2017 workshop on Stranded Assets and Climate Policy; EAERE 2017 (Athens); OxCarre; LSE; BEEER 2017 (Bergen); INFORMS 2017 (Houston); CESifo annual Energy & Climate Economics conference (2017); Misum (Stockholm), envecon2018; SURED (2018) and the 2nd MMCN conference (Stanford) for helpful comments and suggestions. When this project started, Baldwin and Kuralbayeva were employed by the Grantham Research Institute of the London School of Economics and were supported by the UK’s Economic and Social Research Council (ESRC) and the Grantham Foundation for the Protection of the Environment, and Cai was employed by the Becker Friedman Institute of the University of Chicago and a visiting fellow at Hoover Institution at Stanford University. Cai acknowledges support from the National Science Foundation (SES-0951576 and SES-1463644) under the auspices of the RDCEP project at the University of Chicago. Kuralbayeva also acknowledges support from Statoil via the Statoil Chair in Economics at NHH. The authors have no other relevant or material financial interests that relate to the research described in this paper.

<sup>†</sup>Department of Economics and Hertford College, Oxford University, UK; elizabeth.baldwin@economics.ox.ac.uk

<sup>‡</sup>Department of Agricultural, Environmental and Development Economics, The Ohio State University, USA; cai.619@osu.edu

<sup>§</sup>Department of Political Economy, King’s College London, UK; karlygash.kuralbayeva@kcl.ac.uk

# 1 Introduction

When should we make a “green” investment? Is early deployment of expensive options worthwhile, if it helps us to “learn”? And how soon should we ensure that all new investments are “green”? These questions are pertinent across many sectors, especially when investments are long-lived. For example, cities may need to expand to house 1.5 billion more people over the next two decades; characteristics of their current development may affect emissions for many years to come (Bhattacharya et al. 2015). In this paper we focus our attention on one critical source of greenhouse-gas emissions: the power sector.

So, in particular, what is the optimal time to stop investment in fossil-fuel-based power plants? The world continues to make big investments into their construction, particularly coal-based plants: estimates suggest that almost 1 trillion US dollars of such investments are planned (Shearer et al. 2016). Given the long lifetimes of fossil fuel based power plants, the emissions embodied in this infrastructure potentially undermine stringent long-term climate objectives, such as the 2°C target (see Pfeiffer et al. 2016). As such, a fast coal phase-out strategy is considered as one of the necessary conditions to achieve a transformation in line with the Paris Agreement. Some countries such as the UK, Finland and France have significantly reduced their power production from coal in recent years and announced the phasing out of coal completely in the coming 10-15 years. On the other hand, production of electricity from renewable sources has become more competitive, expanding dramatically, primarily due to the decline in costs driven by learning-by-doing in this sector.

Economists advocate a withdrawal from polluting sectors driven by carbon pricing. And investments in the clean energy sector are driven by subsidies. So a natural third question is to ask: which of these policy instruments (carbon tax or subsidy) is more efficient in terms of maximizing social welfare in a second-best setting, where only one instrument is available?

In this paper, we study these questions both theoretically and numerically. Our analysis is thus in two complementary parts. First, in a simplified and very general model, we explore the properties of irreversible investment decisions (Arrow 1968, Arrow and Kurz 1970, Greenwood et al. 1997). The model characterizes optimal irreversible investment decisions when it is known that returns on this capital are due to fall, leading to important implications even without uncertainty. And similarly, we explore investments, returns and optimal subsidies for investments into technologies in a sector that undergoes learning-by-doing (Wright 1936, Arrow 1962).

We then quantify the importance of irreversibility and learning-by-doing in a dynamic general equilibrium climate-economy model. This is based on DICE-2013 (Nordhaus 2014a) and its annualized version (Cai et al., 2016) but deviates in two important ways. Firstly, the energy sector is modeled explicitly, with irreversibility in investments in capital stocks in both a “dirty” sector and a “clean” sector, and with the latter sector characterized by learning-by-doing. And secondly, as well as using the damage function of Nordhaus (2014a), we also consider scenarios in which global temperature changes do not exceed 2°C. This stringent target makes both the irreversibility and the learning-by-doing more important, and it is more in line with current international aspirations. Given the two externalities present in our model (global warming and learning-by-doing), we consider cases in which both a carbon tax and subsidy instruments (the first-best setting) or only one of the two instruments (the second-best) are available.

The four main findings of the paper are as follows. First, we establish a theoretical result for the relationship between climate policies and investment in dirty capital stocks, which we call the “irreversibility effect”: if dirty capital cannot be converted to other capital, then it is optimal to stop investing in dirty capital earlier (compared to a case in which investment is reversible). Irreversibility in investment implies an earlier shift to investment in the clean sector, to avoid a future stranding of assets in the dirty energy sector. This shift therefore reduces emissions in the short term. We

thus demonstrate that irreversibility effects on the demand side *enhance* the effects of a carbon tax in the short term. This is in contrast with the standard Green Paradox effect (see e.g., Sinn 2008, 2015, Jensen et al. 2015) which focuses on the *suppliers* of a fossil fuel resource and shows that the knowledge of an increasing carbon tax will increase extraction of fossil fuels and will thus *counteract* the effects of the carbon tax in the short-term. Moreover, at the time at which investment in dirty fossil fuel infrastructure stops, returns on its existing assets go above those of the general economy. From the perspective of an investor, this result makes perfect sense. In the long-term, returns on this investment will fall, and thus, current investments are only attractive when short-term excess returns (relative to those of the general economy) are sufficient to compensate for future losses.<sup>1</sup> It is natural to ask for how long this “short term” irreversibility effect is relevant. We obtain a lower bound for the time period by which disinvestment must be advanced. If other things are equal, this lower bound is higher for longer-lived assets.

Second, we provide a simple and elegant expression for the optimal subsidy on technologies whose price evolves as a result of “learning-by-doing” occurring in a sector. This subsidy increases with the learning rate, and also with the optimal rate of growth of capital in this sector. (Indeed, it is simply the product of the learning rate with the sum of optimal growth and depreciation). We call this the “*acceleration effect*” for technology policy. Thus suppose, for example, that a carbon tax on substitute sectors makes the clean technology more competitive, and so enhances growth in this sector. The acceleration effect says that the *optimal* growth rate is now faster still: higher growth due to market forces implies a higher optimal subsidy, supporting this high optimal growth. So, the importance of learning-by-doing is accentuated by the early withdrawal from the dirty energy sector.

Third, our quantitative results support our theoretical findings and illustrate that the net (of depreciation) rate of returns on dirty capital infrastructure with irreversible investments follows an unusual trajectory: while initially matching the returns in the general economy, it rises above the returns in the general economy when investment in dirty capital stops, and remains higher for some period of time. Within this period and for some time thereafter, investment in dirty capital will be equal to zero, although the dirty capital is not underutilized. However, returns on dirty capital will fall eventually, reaching zero once this capital is indeed underutilized (so net returns will be negative). Quantitative results illustrate that the timing of these effects depends on the climate policy target: the irreversibility effect is present only if policy objectives are stringent enough.

We can identify the acceleration effect on the optimal subsidy when we compare mild to more stringent targets. With more stringent targets, the renewable sector grows faster. Our theoretical result predicts this, for two reasons: in the stringent case, competing sectors face a higher carbon price, and additionally, the optimal subsidy is higher. We also can observe the acceleration effect by varying the learning rate and calculating the associated optimal subsidy. In early periods, the subsidy is convex in the rate of learning, because technologies whose costs decline faster will also have higher initial optimal growth. However, in later periods (for instance, 2050 or 2100), the relationship is slightly concave. Rapid early deployment of fast-learning technologies leads to faster market saturation, and so slower optimal growth in later periods. Thus, at such times, the subsidy is sub-linear in the learning rate.

Finally, under the second-best setting, we quantitatively explore which instrument—carbon tax or subsidy—yields a lower welfare loss compared with the first-best situation. We show that under a less ambitious climate policy, the economy is better off with the subsidy, while carbon pricing induces a lower welfare loss compared with the subsidy if climate policies are more ambitious.

---

<sup>1</sup>This extra premium on irreversible investment even without uncertainty is also called the commitment premium, see for example Bernstein and Mamuneas (2007).

In terms of the implications of our results, this paper contributes to the debate on characteristics of optimal policy to combat climate change. Some advocate a “gradual slope” in policy implementation because economic growth implies that the current generation is poor relative to future generations, and so should not bear large costs of emission reductions. Moreover, doing so reduces pressure for premature retirement of the existing dirty capital stock, and it provides valuable time to develop low-cost, low-carbon-emitting technologies.<sup>2</sup> Others counter this line of reasoning by arguing that an effective way to reduce abatement costs is to accelerate learning-by-doing.<sup>3</sup> We find that early investment in the renewable sector is crucial, not only to accelerate the decline in the costs of clean energy, but also to prevent later stranding of assets that use fossil fuels. Our quantitative results within the second-best setting (with only one policy instrument available) emphasize the importance of adopting carbon pricing, an instrument that can facilitate a rapid decarbonization of the global power sector under an ambitious climate policy target like that set under the Paris Agreement. However, considering the past 10-20 years, relatively unambitious policy has manifested itself in a large part through subsidies on renewables; if that less ambitious level of emission reductions had been optimal, that single choice of instrument may well have been an excellent second best. Our quantitative results under a second-best setting, where carbon pricing action is restricted by political economy issues, illustrate that the subsidy needs to be higher (if the carbon tax is lower than the optimal first-best level) to compensate. In later periods, once renewable technologies reach maturity, there is no need to maintain the same level of subsidies. Thus, the subsidies help to overcome political constraints in the short-run, to reach the climate target, and help to develop the renewable sector to play its role in the economy in the longer term.

Finally, our paper speaks to the debate on stranded assets and climate policy.<sup>4</sup> The literature so far has focused on the necessity of stranding fossil fuel reserves, if one is to limit climate change to less than two degrees of warming.<sup>5</sup> As we show, the implications for investment are different when one considers the stranding of assets that *use* the fuel.

The rest of the paper is organized as follows. In the next section we discuss related literature. In Section 3 we present a simple analytical model in which we characterize optimal irreversible investment decisions when the anticipated returns on those investments will fall in the future. In Section 4 we consider a simple model of investment into production processes that undergo learning-by-doing. Section 5 describes how we set up the full dynamic general equilibrium climate-economy model to quantify the theoretical results. Section 6 sets out the results from the simulations of the climate-economy model. The final section concludes. Details on the calibration, and proofs of technical results, are provided in the Appendices.

## 2 Related Literature

Our paper is related and contributes to several strands of research. First, irreversibility of investment features prominently in the modern theory of firm-level investment under uncertainty (e.g., Abel 1983, Pindyck 1991, Dixit 1992), and in the “putty-clay” framework of Atkeson and Kehoe (1999).

---

<sup>2</sup>W. Nordhaus was one of those in the past who recommended a “gradual slope”, but he recently argued that a target with a limit of 2°C “appears to be unfeasible with reasonably accessible technologies” (Nordhaus 2018).

<sup>3</sup>Still, some authors find that learning-by-doing has an ambiguous impact on the timing of emissions abatement (Tol 1999, Goulder and Mathai 2000).

<sup>4</sup>Caldecott (2017) argues that the term “stranded assets” has been used to describe various situations. We follow Caldecott et al. (2013) and define stranded assets as assets that have suffered from premature write-down before the end of their technological life.

<sup>5</sup>For instance, according to McGlade and Ekins (2015), an estimated one-third of oil reserves, half of gas reserves and more than 80% of known coal reserves are referred to as “stranded”.

Our work builds on the earlier studies of investment irreversibility in a deterministic setting. The closest precedent to our work is Arrow (1968), who showed that optimal irreversible investment is characterized by alternating periods of positive gross investment and zero gross investment. Our findings are in line with his, but relative to his work, our study applies the irreversibility effect to the case of a polluting industry and demonstrates how irreversibility can affect the path of emissions.

Second, we contribute to the literature on investment under learning-by-doing, which originated with Arrow (1962) and Spence (1981). Significant subsequent studies (Fudenberg and Tirole, 1983; Petrakis et al., 1997) assume that this learning is private; in particular Petrakis et al. (1997) show that in this case all benefits from learning are internalized. Contrastingly, we focus on the pure spillovers case, in which all firms are small and their costs are determined by the cumulative production of the entire sector. Reichenbach and Requate (2012) distinguish between private learning, and that with spillovers, and show that it is learning spillovers that matter. Other parts of this literature also recognize that induced technological change affects optimal climate policy and optimal policy mix between carbon taxes and innovation subsidies (see, e.g., Goulder and Mathai 2000). Kverndokk and Rosendahl (2007) find that optimal renewable subsidies are increasing in the rate of future renewable deployment. However, their findings are numerical results. In contrast, our analysis provides an explicit formula for renewable subsidies derived from a tractable model, thus giving a concrete rationale for how large the renewable sector subsidies should be, and how they evolve over time.

Third, there is an extensive literature that has explored the robustness of the Green Paradox effect by considering various extensions of its typical underlying resource model (e.g., Gerlagh 2011, Michielsen 2014, Smulders et al. 2012, van der Ploeg 2013, van der Ploeg and Withagen 2014). The irreversibility effect, discussed in this paper, complements other mechanisms countervailing against the Green Paradox, discussed in this literature, but it adds a different perspective, as it focuses on the *demand* side of a fossil fuel resource. The only other work exploring this perspective is contemporaneous work by Bauer et al. (2018), who provide a numerical comparison of the irreversibility effect and Green Paradox; we put the phenomenon on a clear theoretical footing.

Fourth, our paper is related and contributes to a large quantitative literature that investigates relative merits of carbon taxes and renewable subsidies to address climate change (e.g., Fischer and Newell 2008, Fischer et al. 2017, Gerlagh and van der Zwaan 2006). However, these studies generally abstract from consideration of different climate policy targets under second-best settings with irreversible investment decisions, as we do in this paper. On the other hand, a rich and growing literature has developed integrated assessment models, to study a number of different climate change issues. Papers assessing future emissions from the energy sector include Pfeiffer et al. (2016) and Davis et al. (2010). However, these do not use the dynamically optimizing frameworks of the economics literature. Other climate-economy models generally ignore the interplay between irreversible investment decisions, inertia in energy systems, and climate policies, on which this paper focuses.

An exception is the concurrent work by Rozenberg et al. (2018), who examine the trade-off between efficiency and political feasibility of climate change mitigation policies, in terms of avoiding stranded assets. The “first best” version of their model is similar to our model of irreversible investments (they do not provide a full calibration, or explore learning-by-doing). However, they do not identify that short-term economic returns on dirty capital can go *above* those in the rest of the economy, after investment has ceased. This is a key part of what we call our “irreversibility effect”. Thus, unlike us, they do not exhibit that dirty capital stock may be accrued in the full expectation that these assets will become stranded. Another difference is that they show that investment must stop (at least temporarily) at the moment at which policy is implemented. This follows from their use of continuous time: with a discrete time step this could simply correspond to reduced investment

within a single period. Our model, calibrated against real-world values, shows positive investment in the short term even under stringent policy scenarios.

Finally, our paper belongs to the literature on path dependence and climate change (Fouquet 2016; Aghion et al. 2014, 2016; Grubb et al. 1995; Vogt-Schilb et al. 2018). We contribute to this literature by analyzing the implications of path dependence embodied in carbon-intensive infrastructure for the design of optimal climate change policies under first- and second-best settings.

### 3 A Simple Model of Irreversible Investments

As described in the introduction, our key motivation is studying capital stock effects in the context of climate change. However, the analytical results we prove in this context hold in much more general settings. Therefore, in both this section and the next, we present models that focus only on the “moving parts” which are relevant to these results. The model analyzed here, and the model of learning-by-doing analyzed in Section 4, will both be embedded in our full structure in Section 5.

For this section, we consider the implications of irreversibility, specifically in investments in capital stocks whose economic productivity will decline (cf. Arrow 1968, Arrow and Kurz 1970).<sup>6</sup> Our key example is investment in fossil-fuel-using power stations, but as discussed earlier, many other illustrations exist, such as the structure of transportation systems and the design of cities.

#### 3.1 The Household’s problem

Consider a representative household, which holds  $k_t$  of a certain “irreversible investment” asset. The household can make an additional investment of  $i_t \geq 0$  in each period  $t$ . We will assume that the period- $t$  return  $r_t$  on the asset eventually drops below the returns in the general economy and we ask when investment stops, i.e. from which point in time we have  $i_t = 0$ .

There are other opportunities for investment and other sources of income, written net as  $o_t$ . The household’s per-period consumption is  $c_t$ . Their budget constraint is  $i_t + c_t = r_t k_t + o_t$  where investment  $i_t = k_{t+1} - (1 - \delta)k_t$ , with  $\delta$  being depreciation of the irreversible investment asset. Write  $u$  for their utility function and  $\beta$  for their utility discount factor. We will want to compare returns on the irreversible investment asset with those on the rest of the economy. To do this without specifying any more details, just write the standard ratio from the Euler equation as  $e_{t+1} := \frac{u'(c_t)}{\beta u'(c_{t+1})} - 1$ , which represents the consumption discount rate. Make the minor assumption that  $e_t$  is bounded, and bounded away from  $-\delta$ . That is, assume there exist  $\epsilon > 0$  and  $R \gg 0$  with  $-\delta + \epsilon < e_t < R$  for all  $t$ .

We first see that the Euler equation does indeed hold while investment is positive. See Appendix A for proofs of all results in this section.<sup>7</sup>

<sup>6</sup>Our irreversibility is an example of a “putty-clay” framework. This framework assumes that dirty-electricity-producing firms can turn variable raw capital that is malleable ex ante (“putty”), into capital goods with certain technological characteristics, including energy efficiency. However, once the choices have been made, and the coal-fired plant is in operation, there is a fixed factor ratio, or frozen structure (“clay”). This property of capital structures implies that the plant will either be fully utilized or scrapped at some point in time when it loses its economic value. The relevant literature that uses putty-clay production functions are Atkeson and Kehoe (1999) and Casey (2017). The putty-clay framework should allow simultaneous use of old and new vintages and is particularly useful in explaining phenomena like entry and exit of firms, especially through the fact that equipment has become obsolete in an economic sense, or to explain drives of changes in final-use energy intensity in the US as done by Casey (2017).

<sup>7</sup>Lemmas 3.1-3.3 follow straightforwardly from a technical lemma (Lemma A.1 in Appendix A) that presents the shadow price on the irreversibility constraint,  $i_t \geq 0$ , as the net present value of investment in this asset, relative to the opportunity cost.

**Lemma 3.1.** *For any times  $s_0$  and  $s_1 > s_0$ , investment  $i_t > 0$  holds for all  $t \in \{s_0, \dots, s_1\}$  only if  $r_t - \delta = e_t$  for  $t \in \{s_0 + 1, \dots, s_1\}$ .*

But now let us consider a more interesting case. Suppose the net return from the irreversible investment asset drops *below*  $e_t$  at some time  $s_1$ , and this is anticipated. Changing economic conditions mean that this asset is no longer as productive as it was. Then we stop investing *earlier* than time  $s_1$ , and reap *excess* returns on the irreversible investment for some of the intervening period.

**Lemma 3.2.** *Suppose that initial investment  $i_0 > 0$  and that it is correctly anticipated that  $r_t - \delta < e_t$  for  $t \in \{s_1, \dots, s_2\}$ . Then there exists a point  $s_0$  in time such that  $s_0 \leq s_1 - 1$ , such that  $r_{s_0} - \delta > e_{s_0}$  and such that investment  $i_t = 0$  for  $t \in \{s_0, \dots, s_2 - 1\}$ .*

Lemma 3.2 tells us that net returns  $r_t - \delta$  from this asset follow an unusual trajectory: initially matching the consumption discount rate path  $e_t$ , we see that  $r_t - \delta$  rises above  $e_t$  at some point before it falls beneath. Investment is zero while returns follow this pattern. (This pattern is illustrated in Section 6, Figure 2).

From the perspective of an investor, these short-term excess returns make perfect sense. If the investor knows that, in the long-term, returns on this infrastructure will fall, then it is not an attractive investment. However, short-term additional gains will compensate for long-term losses. In fact, if we write  $\Delta_{t,s} = \prod_{s'=1}^s \frac{1}{1+e_{t+s'}}$  for the compound consumption discount factor, then it must hold that:

**Lemma 3.3.** *Suppose that initial investment  $i_0 > 0$ . Then for any times  $s_1$  and  $s_2$  (e.g., those as in Lemma 3.2), it holds that*

$$\sum_{s=1}^{s_1-1} (1-\delta)^{s-1} \Delta_{0,s}((r_s - \delta) - e_s) \geq \sum_{s=s_1}^{s_2} (1-\delta)^{s-1} \Delta_{0,s}(e_s - (r_s - \delta)) \quad (1)$$

Moreover, these short-term gains will indeed be realized if all other investors are similarly ending investment early.<sup>8</sup>

How early does investment stop? Lemma 3.2 could allow  $s_0 = s_1 - 1$ . But Lemma 3.3 shows that this is unlikely: we need short-term gains to compensate the long-term losses. It is hard to find a closed form expression for the time gap in general. But if we put bounds on  $r_t - \delta$ , and if we assume that  $e_t$  is a constant  $e$  so that  $\Delta_{0,s}$  is just  $\frac{1}{(1+e)^s}$ , we obtain a bound for the advancement of the end to investment:

**Corollary 3.4.** *Suppose  $e_t = e$  is constant. If initial investment  $i_0 > 0$ , if net returns  $r_t - \delta \leq e - d_2$  for all times  $t \geq s_1$ , and net returns have an overall bound, that is, there exists  $d_1 > 0$  such that  $(r_t - \delta) \leq e + d_1$  for all  $t$ , then there exists a point  $s_0$  in time such that  $s_0 \leq s_1 - 1$ , such that investment  $i_t = 0$  for all  $t \geq s_0$  and such that*

$$s_1 - s_0 \geq \frac{\log(d_1 + d_2) - \log(d_1)}{\log(1+e) - \log(1-\delta)}. \quad (2)$$

---

<sup>8</sup>Our results can be illustrated with a historical example, for which we are grateful to Roger Fouquet. In the first half of the 19th century, the introduction of steam engines brought cheaper and more comfortable medium and longer distance travel than had previously been provided by stagecoaches (pulled by horses). Coach companies responded to this heightened competition from railways by ceasing investment into equipment and horses, driving their prices even higher, which inevitably accelerated the transition to railways (Fouquet, 2012).



Note that (2) only provides a lower bound on the advancement of the end of investment. Nonetheless, it provides useful verification of natural intuitions. The right hand side of (2) increases with  $d_2$  and decreases with  $d_1$ . That is, if future returns from the irreversible asset will be extremely low, or if short-term gains are very limited, then this advances the date by which investment must have stopped. Moreover, the right hand side of (2) is decreasing with  $\delta$ . So, *ceteris paribus*, our time by which investment must end is sooner for longer-lived assets.

Thus, for example, suppose the time step is one year. If the Euler rate  $e$  is 0.1, if depreciation is fast at 0.1, and if short-term gains can exceed the Euler rate by the same amount as the deficit in long-term losses, then  $s_1 - s_0$  is bounded below by only 4. But if  $e = 0.04$ , if depreciation is a much slower at 0.01 (considering perhaps the design of cities), and if short-term gains can only achieve one third of long-term losses relative to  $e$ , then investment must end at least 29 years before these low returns commence.

Now we consider what this means for the quantity of total holdings of this irreversible asset. If there are  $L_0$  identical households in the economy, each of size  $\frac{L_t}{L_0}$  at time  $t$ , then the investment behavior of Lemma 3.2 simply scales up. We use capital letters to denote total capital  $K_t$  and total investment  $I_t$  for the irreversible asset. We assume, as is implied by standard models, that in each period, returns  $r_t$  are monotone strictly decreasing in capital stock  $K_t$ . So the pattern of investment implied by Lemma 3.2 implies *a short-term decrease in the dirty energy capital stock, relative to a world in which investments are reversible (and so underutilization is never an issue)*.

To explore this, consider an otherwise identical model in which we relax the constraint  $i_t \geq 0$ —allowing holdings of this capital stock to be converted back into cash for consumption or other purposes. We use tilde to refer to variables in this modified model ( $\tilde{K}_t$ ,  $\tilde{I}_t$ , and so on). We suppose that  $e_t$  is unchanged by relaxing the constraint  $i_t \geq 0$ , because the sector concerning the irreversible asset is very small in relation to the rest of the economy.

**Corollary 3.5.** *Suppose that total initial investment  $I_0 > 0$  and suppose that there exists a point  $t_1$  in time such that  $t_1 \geq 1$  and such that  $\tilde{I}_{t_1} < 0$ . Then there exists a point  $t_0$  in time such that  $t_0 < t_1$ , such that  $I_{t_0} < \tilde{I}_{t_0}$ , and such that  $K_t < \tilde{K}_t$  for  $t \in \{t_0 + 1, \dots, t_1\}$ .*

That is, in the short term, less is invested in the irreversible capital stock, relative to a world in which investments are reversible. By making the same assumptions as in Corollary 3.4 about net returns in the irreversible world, we can provide the same bound on the length of time for which investments are lower than they would be if we could assume reversibility.

### 3.2 The Irreversibility Effect in Climate Change Economics

In this paper we apply the observations of Section 3.1 to a model of climate change economics. We are particularly concerned with capital investments in installations, such as coal fired power stations, which will burn fossil fuels. The quantity of fuel demanded, and burnt, is associated with the quantity of appropriate capital infrastructure available and in use. If, in the extreme case, this relationship is Leontief, then Corollary 3.5 implies:

**Corollary 3.6. [The Irreversibility Effect]** *Suppose emissions are directly proportional to the utilized fraction  $\zeta_t$  of installed infrastructure that uses fossil fuel. Assume that investment in this infrastructure is non-zero in the first period, but there exists a point  $t_1$  in time such that  $t_1 \geq 1$ , and such that this infrastructure would be globally divested if it could be. Then, for some period leading up to  $t_1$ , emissions are below the level they would reach if divestment were possible.*

That is, capital stock effects within those who demand fossil fuels enhance the effect of a carbon tax in the short term. Again, we can turn to Corollary 3.4 to give a bound on the length of this “short term”.

This result contrasts with the Green Paradox (Sinn, 2008), relating to capital stocks in the *supply* of fossil fuels. Suppose a new carbon tax regime has just been announced, in which future carbon taxes are higher than had previously been expected. This reduces the future rents from stocks of fossil fuels. So fossil fuel suppliers update their optimal extraction pathways, reducing their short term per-unit rents to increase the volume sold. Short-term emissions are, thus, higher than they would have been estimated by a more naive model, which ignored these supply-side effects (and a Green Paradox occurs when this effect is large enough to increase short-term emissions). However, as developed above, investment irreversibilities mean that short-term emissions are lower than they would have been estimated by a more naive model, which ignored these demand-side effects. Thus these two effects push in countervailing directions, both in a richer model and in the real world.

We emphasize the contrast theoretically in this paper. Our numerical methods (Sections 5–6) are not well-tuned to assess which of these effects is more important, because our stylized model only distinguishes one fossil fuel, which is calibrated to coal (see Section 5.5). However, concurrent complementary research (Bauer et al. 2018) has examined this question numerically. There, models with much more finely specified energy sectors show that the Green Paradox primarily applies to oil, rather than coal (although a small Green Paradox can exist for coal); but that the irreversibility effect, which is more important to coal, is also generally much greater. Their only scenario in which the Green Paradox dominates the irreversibility effect has a carbon tax which is very low, with a 40 year implementation lag.

Our main interest, as in Corollary 3.6, is in the effect of irreversibilities on the timing of emissions (Figure 3). Corollary 3.6 does not address how Pigouvian taxes change, once irreversibilities are taken into account. A simplistic understanding would suggest that carbon taxes may be lower, since emissions are reduced in the short term, by the irreversibility effect. However, the question is not this straightforward. The proof of Corollary 3.5 showed that incorporating irreversibility also means that there will be greater holdings of that asset at some date after the time  $t_1$  at which infrastructure would be divested if this were possible. It does not necessarily follow that emissions will be greater after this date, because we allow underutilization of capital stocks. (For the same reason, the long-term implications of the irreversibility effect are ambiguous.) But if emissions are indeed higher after date  $t_1$  in the irreversible scenario, and if there is a reasonably low discount rate, then the net effect on the Pigouvian tax is reduced. We have found in the calibration of Section 5 that the difference between an optimal carbon tax with and without irreversibility is small.

## 4 A Simple Model of Investing with Learning-By-Doing

Learning-by-doing is often cited as a rationale for subsidizing renewable electricity. The theory of learning-by-doing is motivated by simple observation: production performance (either in the form of productivity or cost reduction of technology) tends to improve with the accumulation of experience. We are particularly interested in the form that was specified by both Wright (1936) and Arrow (1962): each doubling of cumulative deployment reduces costs (and hence prices) by the same factor, known as the “learning rate”.<sup>9</sup> Empirically, the existing literature has found substantial evidence that the price of renewable energy evolves in this way, although a causal relationship has not been finally established.<sup>10</sup> For this paper we assume that this causality does indeed hold. We

<sup>9</sup>Wright (1936) was the first one to describe the concept of learning, after observing a uniform decrease in the number of direct labor hours required to produce an airframe for each doubling of the cumulative production of the plant under consideration.

<sup>10</sup>Lindman and Soderholm (2012) use aggregate data and show that learning externalities are present in wind turbines and solar panel costs. Such studies based on aggregate data, however, are unable to disentangle the effect of exogenous technological change from the effect of leaning-by-doing, thus masking the diverse drivers of technology

also assume that firms are small enough so that they do not internalize learning-by-doing effects when making decisions.

To model this, we consider stocks of a ‘learning-by-doing (LBD) asset’,  $H_t$  (writing  $h_t$  for household-level holdings as before). The notation reminds readers that the LBD asset embodies human capital in the form of knowledge, as well as the infrastructure itself. The form of this knowledge is embodied in the price  $p_t^H$  of installing this infrastructure:  $i_t^H = p_t^H(h_{t+1} - (1 - \delta)h_t)$ . This price depends on the total installed capacity  $H_t$ , which aggregates individual holdings  $h_t$ . Thus,  $p_t^H = G(H_t)$ . Of particular interest is Wright’s Law: there exists a constant  $\lambda > 0$  with

$$p_t^H = G(H_t) = p_0^H \left( \frac{H_t}{H_0} \right)^{-\lambda}. \quad (3)$$

However, we also derive the optimal subsidy for the general case, minimizing additional assumptions. Investment in this sector is, again, irreversible.

Learning-by-doing gives rise to an externality. So we first explore the optimal program of investment found by a social planner. We contrast this with the behavior of households who act as price-takers, to identify and understand the optimal subsidy.

#### 4.1 The Social Planner’s Case

The social planner maximizes total welfare  $\sum_{t=0}^{\infty} \beta^t L_t u \left( \frac{C_t}{L_t} \right)$ , where  $L_t$  is the population size and  $C_t$  is total consumption. This is subject to the investment equations  $I_t^H = p_t^H(H_{t+1} - (1 - \delta)H_t)$ ; the investment bounds  $I_t^H \geq 0$ ; the price evolution  $p_t^H = G(H_t)$ ; and the budget constraints  $I_t^H + C_t = f_t(H_t, O_t)$ , where we have written  $O_t = L_t o_t$  for the economy-level aggregate of “other” incomes, so that we can write  $f_t(H_t, O_t)$  for the production function. The planner will treat  $O_t$  as exogenous, which is a harmless assumption if all externalities in the remainder of the economy have been internalized.

Define the *direct return* on investments in the LBD asset to be  $r_{t+1}^s := \frac{1}{p_{t+1}^H} \frac{\partial}{\partial H_{t+1}} f_{t+1}(H_{t+1}, O_{t+1})$ .

That is, we account for the price of investments.<sup>11</sup> Because this investment price changes over time, the direct return does not reflect the full value to the social planner of investments in our LBD asset. So we use the discrete time version of the definition of Jorgenson (1967) to define a *shadow return* on these investments:

$$R_{t+1} := \frac{\mu_t^H - \beta(1 - \delta)\mu_{t+1}^H}{\beta u'(C_{t+1}/L_{t+1})}$$

Here,  $\mu_t^H$  is the shadow price on the investment equation  $I_t^H = p_t^H(H_{t+1} - (1 - \delta)H_t)$ , and so gives the total shadow value of marginal investments in the LBD asset. To find the return realized in period  $t + 1$ , we subtract the discounted depreciated shadow value going further forward. As usual, everything is measured relative to the marginal value today of consumption tomorrow.

To show how natural the shadow return is, and relate it to the direct return, we write again  $e_{t+1} := \frac{u'(C_t/L_t)}{\beta u'(C_{t+1}/L_{t+1})} - 1$  and show:<sup>12</sup>

---

costs (see also Nordhaus 2014b). Nemet (2006) for instance finds that after accounting for measures of technological change and the cost of inputs, learning has only weak explanatory power for solar panel costs. Much more recently, Lafond et al. (2018) use hindcasting techniques to assess this model, and find that it provides a very good fit. Bollinger and Gillingham (2014) provide evidence for cost reductions due to learning-by-doing across installation contractors of solar photovoltaics in California from 2002 to 2012. See also Rubin et al. (2015) on the use of learning rates in this context.

<sup>11</sup>We write  $r_t^s$  to distinguish from the notation for the market rate of return  $r_t$ , which we will use in Section 4.2.

<sup>12</sup>See Appendix A for proofs of all results in this section.

**Proposition 4.1.** *Suppose that investment into this sector will be non-zero next period, i.e.  $I_{t+1}^H > 0$ . Then  $R_{t+1} - \delta = e_{t+1}$ , and*

$$\underbrace{\frac{p_t^H}{p_{t+1}^H} R_{t+1}}_{\text{shadow return}} = \underbrace{r_{t+1}^s}_{\text{direct return}} - \underbrace{\frac{p_t^H - p_{t+1}^H}{p_{t+1}^H} (1 - \delta)}_{\text{price effect}} - \underbrace{(H_{t+2} - (1 - \delta)H_{t+1}) \frac{G'(H_{t+1})}{p_{t+1}^H}}_{\text{learning effect}} \quad (4)$$

Thus, the Euler equation holds for *shadow* returns (as distinct from direct returns, for which it does not hold). The factor  $p_t^H/p_{t+1}^H$  on  $R_{t+1}$  in (4) is needed because, as defined,  $R_{t+1}$  values returns relative to the price  $p_t^H$  of investment at the moment at which the investment is made, while direct returns  $r_{t+1}^s$  are valued relative to the price  $p_{t+1}^H$  at the time at which we receive the return.

Relative to next-period prices, then, the shadow return is composed of three terms. One is the “direct return”. Next, we observe a “price effect”, from the dependence of prices on time. In period  $t + 1$ , an additional unit of renewable capital costs  $p_{t+1}^H$ , but it would have cost  $p_t^H$  in period  $t$ . If prices are decreasing over time then this gives an incentive to delay investment, and so reduces the shadow return on investment in period  $t$ . This effect arises because we assume that prices are constant within each period, and so the benefits of learning are not enjoyed until the following period. Thus, the importance of this consideration will depend on the size of the time-periods we use. In our calibration (Section 5) the time step is one year, which may be reasonable: technological innovations cannot be shared instantaneously, but realized prices do seem to differ from year to year.

The final term, which we call the “learning effect”, arises due to our assumption of learning-by-doing. It incorporates the marginal change in price in the LBD asset due to our holdings of this asset, valued against their price at time  $t + 1$ . This marginal change in price is multiplied by how many units of the asset we will invest in, in period  $t + 1$ . One must not be confused by the negative sign: typically  $G'(H) < 0$  (prices decrease with capacity), and  $H_{t+1} > (1 - \delta)H_t$  (investment is positive), so that the learning effect is typically positive.

The net effect of the price and learning effects may be positive or negative, and so the total return on renewables may be greater than, or less than, their direct net return (see Corollary A.2 for examples of each case and a discussion). However, the price effect will be taken into account by small rationally optimizing firms, whereas the learning effect will not, because in our specification, learning-by-doing is a pure externality. So, as we will see next, the optimal subsidy in a decentralized model is equal to the learning effect. It follows that investing in the LBD asset becomes worthwhile from a social perspective before it is individually rational: if investment will take place in the near future, it is socially optimal to start earlier than an individual would choose to.

## 4.2 Learning-By-Doing and the Acceleration Effect

We assume that households act as price-takers on the LBD asset. Due to the positive externality, there will be under-investment without intervention. So we introduce a subsidy,  $\tau_t$ ; it is convenient to express this as a subsidy on the rate of return. Now we may write the household’s budget constraint as  $i_t^H + c_t = (r_t + \tau_t)p_t^H h_t + o_t$ , where  $o_t$  represents other sources of income (as in Section 3.1). These investments are characterized by  $i_t^H = p_t^H(h_{t+1} - (1 - \delta)h_t)$  and  $i_t^H \geq 0$ . The subsidy is financed by lump sum taxation; as the households are price-takers, this taxation may be incorporated into  $o_t$ . Meanwhile, a final goods firm maximizes its profits  $f_t(H_t, O_t) - r_t p_t^H H_t - p_t^O O_t$ , where  $p_t^O$  is the price they must pay for access to other assets.

Again there are  $L_0$  households in the economy, each of size  $\frac{L_t}{L_0}$  at time  $t$ , so that the consumption of a representative individual is  $\frac{L_0 c_t}{L_t}$ . Write  $g_t^H$  for growth  $\frac{H_{t+1} - H_t}{H_t}$ . Then:

**Proposition 4.2.** *Suppose that any externalities in  $o_t$  have been internalized. The subsidy  $\tau_t$  which maximizes consumer welfare  $\sum_{t=0}^{\infty} \beta^t L_t u\left(\frac{L_0}{L_t} c_t\right)$  is equal to the learning effect under the optimal growth path:*

$$\tau_t = -(H_{t+1} - (1 - \delta)H_t) \frac{G'(H_t)}{p_t^H} = -(g_t^H + \delta) \frac{H_t}{p_t^H} G'(H_t)$$

This expression is even simpler if  $G(H_t)$  follows Wright’s Law (3).

**Corollary 4.3.** [*The Acceleration Effect*] *If  $G(H_t) = p_0^H \left(\frac{H_t}{H_0}\right)^{-\lambda}$ , then*

$$\tau_t = \lambda (g_t^H + \delta).$$

*In particular, the optimal subsidy  $\tau_t$  increases with optimal growth  $g_t^H$ .*

Thus, the subsidy to the LBD asset is a straightforward function of its learning rate and its optimal growth rate. Contrary to models which prescribe a short-term subsidy to this sector, the specification we use implies that this subsidy is positive as long as there is any investment in this sector, even only to replace depreciating stock.

Moreover, if a change in information or policy makes the LBD asset more attractive in the economy, and so it starts to accumulate faster, *irrespective* of the subsidy, it follows that the optimal subsidy is higher than it was before. This higher subsidy will *also* increase growth of the sector, and we conclude that the *optimal* growth in the sector has increased by a greater extent than what was induced by the new-found comparative advantage—because that comparative advantage makes learning in this sector even more socially advantageous. We call this the *acceleration effect* for technology policy.

For an illustration, see Figure 4a in Section 5.9, where we consider the optimal subsidy under both ‘mild’ and ‘stringent’ climate policy targets (as will be precisely defined in Section 5.7, below). There we see that, in the short term, more ambitious targets to decarbonize the economy, which will incentivize faster deployment of renewable technologies, also increase the optimal subsidy to investments in these technologies. In the longer term, stringent climate policy targets mean that we have already developed a greater capacity of this capital stock, and so its optimal growth rate drops to a lower level than under mild policy targets; the subsidy is therefore also lower.

In the same way, the optimal initial subsidy is convex in the learning rate: if the learning rate is higher for a technology, then *ceteris paribus* optimal growth in this technology is also higher, and hence the subsidy is higher for this reason as well as because of its straightforward dependence on the learning rate. See Figure 5 for an illustration and further discussion.

## 5 The Full Model

This section outlines the full dynamic general equilibrium climate-economy model which is used for quantitative analysis. The derivations of the equations that define the solution of the model are given in Appendix D. To summarize, the model presents a climate-economy structure, where, unlike many leading climate-economy models, we differentiate between three capital stocks: general capital, “clean”, and “dirty”.<sup>13</sup> Irreversibility in investments characterizes the latter two capital

<sup>13</sup>In a similar way, but within a different context, Greenwood et al. (1997) developed the importance of investment in differentiated capital stocks for growth and technological change. For climate-economy models which do not differentiate capital stocks, see, for example, Nordhaus (2008); Golosov et al. (2014); Rezai and van der Ploeg (2017); Acemoglu et al. (2016); Barrage (Forthcoming); Cai and Lontzek (Forthcoming) and Cai et al. (2018).

stocks, as in Section 3 above. We allow underutilization of dirty capital stocks, once they become uncompetitive. In addition we assume that the “clean” sector is characterized by “learning-by-doing”: costs of new technologies decline as a function of cumulative installed capacity in the sector, as in Section 4. The climate module uses the representation of the carbon cycle, temperature system, and climate-economy feedbacks based on the DICE framework (Nordhaus, 2014a), but calibrated to an annual time step (Cai et al., 2016).

There are five production sectors and, thus, there are five types of firms: final-goods producing firms, aggregate-electricity producing firms, dirty-electricity producing firms, fossil-fuel extracting firms and firms producing electricity from renewable sources. All firms operate under perfect competition.

Turning to the demand side of the economy, we are interested in the behavior of a representative household who does not internalize the learning-by-doing externality and treats all prices as given. Finally, there are three sources of carbon dioxide emissions: general output production, electricity production from dirty energy inputs, and land use. Climate change affects productivity in the sector producing the final goods.

### 5.1 The households’ problem

We are interested in the behavior of a representative household. There are  $L_0$  households (defined as the population size of the economy at the initial period, which in our calibrated model is 2012), and the size of the family at time  $t$  is  $\frac{L_t}{L_0}$ , where  $L_t$  is the population size at period  $t$ .<sup>14</sup>

We consider all variables on a per-household basis, denoted using lowercase letters, while capital letters denote aggregate variables (over all households). For instance, we will write  $k_t^g = \frac{K_t^g}{L_0}$ , where  $K_t^g$  is the aggregate general capital stock and  $H_t$  is the aggregate stock of renewable energy knowledge and capital. The household seeks to maximize the sum of the welfare of individual family members, that is:

$$\sum_{t=0}^{\infty} \beta^t \frac{L_t}{L_0} u\left(\frac{C_t}{L_t}\right) = \sum_{t=0}^{\infty} \beta^t \frac{L_t}{L_0} u\left(\frac{L_0}{L_t} c_t\right)$$

where  $C_t$  is aggregate consumption and  $c_t := \frac{C_t}{L_0}$  is per-household consumption. The household owns a representative share of all three capital assets and the five sorts of companies. We denote  $r_t^D$ ,  $r_t^H$  and  $r_t^g$  as, respectively, the rate of return on capital assets in fossil-fuel-using (dirty) capital; renewable (clean) capital; and general capital (used in the production of final-goods producing firms). Further, we write  $w_t$  for the wage;  $\Pi_t^g$  for the total profit from the sale of the final goods;  $\Pi_t^D$  for the total profit from the sale of fossil-fuel-based electricity;  $\Pi_t^H$  for the total profit from the sale of “clean” electricity;  $\Pi_t^{DE}$  for the total profit from the sale of fossil fuels; and  $\Pi_t^E$  for the total profit from the sale of aggregate electricity. So the aggregate profit is  $\Pi_t = \Pi_t^g + \Pi_t^D + \Pi_t^H + \Pi_t^{DE} + \Pi_t^E$ , and the per-household profit is  $\pi_t := \frac{\Pi_t}{L_0}$ .

In each period, the household faces the following budget constraint:

$$\begin{aligned} i_t^g + i_t^D + i_t^H + c_t &= \frac{L_t}{L_0} w_t + \pi_t + r_t^g k_t^g + r_t^D p_t^D k_t^D + r_t^H p_t^H h_t \\ &+ \frac{1}{L_0} (\tau_t^D (D_t^E + D_t^g) - \tau_t^H p_t^H H_t) \end{aligned}$$

where  $i_t^g$  is investment in general capital,  $i_t^D$  is investment into dirty capital used in the production of dirty electricity,  $i_t^H$  is investment in capital used in the clean sector,  $k_t^g, k_t^D, h_t$  are capital stocks

<sup>14</sup>Table A.2 in the Appendix B provides a summary of variables’ notation and definition.

in the general, dirty and clean sectors respectively,  $\tau_t^D$  is the carbon tax,  $\tau_t^H$  is the subsidy, and  $D_t^E$  and  $D_t^g$  are carbon emissions in the dirty and general sectors respectively. Since we measure fossil and renewable energy capital in gigawatts (GW),  $p_t^D$  and  $p_t^H$  are the respective prices of fossil fuels and renewable energy capital in \$/GW. The price  $p_t^H$  of renewable energy capital falls with our embodied technological progress in the renewable energy knowledge and capital stock, and evolves as (Arrow, 1962):

$$p_t^H = G(H_t) = p_0^H \left( \frac{H_t}{H_0} \right)^{-\lambda} \quad (5)$$

However, as we treat a household as very small, we assume that their investment in renewable energy capital does not influence its price, so that the learning-by-doing externality arises. That is, the household takes  $p_t^H$  as given. The price of fossil fuel capital will be fixed, so that  $p_t^D = p^D$ .

Finally, we assume that the household receives rebates on the taxes and pays for the subsidies (the last two terms in the right hand side of the budget constraint), but as we assume households are small they cannot affect these levels.

The capital stocks in the general, dirty and renewable sectors are accumulated according to the following equations respectively:

$$\begin{aligned} i_t^g &= k_{t+1}^g - (1 - \delta^g)k_t^g \\ i_t^D &= p_t^D(k_{t+1}^D - (1 - \delta^D)k_t^D) \\ i_t^H &= p_t^H(k_{t+1}^H - (1 - \delta^H)k_t^H) \end{aligned}$$

where  $\delta^g$ ,  $\delta^D$ , and  $\delta^H$  are depreciation parameters, and

$$\begin{aligned} i_t^D &\geq 0 \\ i_t^H &\geq 0 \end{aligned} \quad (6)$$

are the irreversibility assumptions - a non-negativity constraint on the rate of accumulation of both dirty and clean capital.

## 5.2 The final-goods firms' problem

The final goods are produced by identical firms, but output is damaged by climate change. Because this sector exhibits constant returns to scale, we can work with aggregate variables, and so write output:

$$Y_t = \Omega(T_t)f(Y_t^g, E_t) \quad (7)$$

where  $T_t$  is the temperature change from pre-industrial levels,  $\Omega(T_t)$  is the damage factor ( $1 - \Omega(T_t)$  is the ratio of damage to output),  $E_t$  is electricity and  $Y_t^g$  is "general" (non-electricity) output.

The final-goods firms individually maximize their discounted profits, so that on aggregate:

$$\sum_{t=0}^{\infty} q_t \Pi_t^g = \sum_{t=0}^{\infty} q_t \left( \Omega(T_t)f(Y_t^g, E_t) - r_t^g K_t^g - w_t L_t - p_t^e E_t - \Psi_t - p_t^{fuel} D_t^g \right)$$

where  $q_t := \beta^t \frac{u'(c_t)}{u'(c_0)}$  is a compound discount factor for the relative price of consumption in period  $t$ , expressed in period 0 units.<sup>15</sup> To produce final goods, these firms rent (aggregate) capital  $K_t^g$ , hire labor  $L_t$ , purchase aggregate electricity  $E_t$  at price  $p_t^e$ , and buy fossil fuel  $D_t^g$  from fossil fuel

<sup>15</sup>See Appendix D.1 for more detailed discussion on the derivation of compound interest for the firms' problems.

extracting firms at price  $p_t^{fuel}$ . The firms spend money on abatement  $\Psi_t$ , which is assumed to abate a fraction  $\eta_t$  of emissions via the following relation:

$$\Psi_t = \frac{\phi_{1,t}\eta_t^{\phi_2}}{(1-\eta_t)^{\phi_3}}Y_t^g$$

so that the emissions constraint is given by:

$$D_t^g = \sigma_t(1-\eta_t)Y_t^g$$

where  $\phi_2$  and  $\phi_3$  are parameters and  $\sigma_t$  represents the ratio of carbon-equivalent emissions to output, all of which evolve exogenously along with the parameter  $\phi_{1,t}$ , as in Cai et al. (2016). Firms do not take into account their emissions' impact on the pollution stock and, thus, on productivity. In other words, firms take  $\Omega(T_t)$  as a given. This, in a conjunction with the knowledge externality in the renewable sector, represents a “twin-market failure” (Jaffe et al., 2005).

For the solution of the model, we assume that the function for production before damages takes the constant elasticity of substitution (CES) form (Hassler et al., 2012):

$$f(Y_t^g, E_t) = \left[ (1-\theta)(Y_t^g)^{1-1/\kappa} + \theta(E_t)^{1-1/\kappa} \right]^{\frac{1}{1-1/\kappa}}.$$

and

$$Y_t^g = f_t^g(K_t^g, L_t) = A_t^g(K_t^g)^\alpha(L_t)^{1-\alpha}.$$

Here,  $\theta$ ,  $\kappa$ , and  $\alpha$  are parameters,  $A_t^g$  is a technology process in the general sector,  $K_t^g$  is general capital and  $L_t$  is labor. Both  $A_t^g$  and  $L_t$  evolve exogenously in the same way as in Cai et al. (2016).

### 5.3 The aggregate-electricity-producing firms' problem

These firms again face constant returns to scale, so we can work with aggregate variables. They produce aggregate electricity  $E_t = f_t^E(H_t, \Gamma_t^{ED})$  which is a combination of fossil fuel production capacity  $\Gamma_t^{ED}$ , and clean production capacity  $H_t$ , with these inputs being priced at  $p_t^{EH}$  and  $p_t^{ED}$  respectively. They sell their output at price  $p_t^e$ , so that the firms maximize the present value of their profits, so that on aggregate:

$$\sum_{t=0}^{\infty} q_t \Pi_t^E = \sum_{t=0}^{\infty} q_t (p_t^e f_t^E(H_t, \Gamma_t^{ED}) - p_t^{EH} H_t - p_t^{ED} \Gamma_t^{ED})$$

In modeling the electricity sector, we follow Papageorgiou et al. (2017)<sup>16</sup> and assume a CES production function of renewable production capacity  $H_t$  and dirty production capacity  $\Gamma_t^{ED}$ :

$$E_t = f_t^E(H_t, \Gamma_t^{ED}) = A_t^E \left( \omega H_t^\xi + (1-\omega)(\Gamma_t^{ED})^\xi \right)^{1/\xi}, \quad (8)$$

where  $A_t^E$  is a technology process in the electricity sector and  $\omega$  and  $\xi$  are CES parameters.

---

<sup>16</sup>We do not use the version of their model in which overall energy is a combination of electricity and “other dirty energy”, as in their model the latter requires no capital input and so is disproportionately favored under optimization.



#### 5.4 The dirty-electricity-producing firms' problem

The dirty electricity producing firms are fossil-fuel-based power stations, which combine existing infrastructure (such as coal-based power plants) with fossil fuels via a Leontief production function. Again, due to constant returns, we may work at the aggregate scale:

$$\Gamma_t^{ED} = \min[\zeta_t K_t^D, D_t^E / \nu]$$

where  $K_t^D$  is total capital in dirty electricity production,  $\zeta_t \in [0, 1]$  is the utilization rate, and  $\nu$  is the conversion rate from fossil fuel to electricity. The Leontief function implies a fixed ratio between utilized fossil fuel energy capital and dirty fuel use:

$$D_t^E = \nu \zeta_t K_t^D. \quad (9)$$

The firms buy fossil fuel  $D_t^E$  at price  $p_t^{fuel}$ , rent the dirty capital infrastructure at rate  $r_t^D$ , and sell their output  $\Gamma_t^{ED}$  to the aggregate electricity producing firms at price  $p_t^{ED}$ . So, the firms in this sector maximize the present value of their profits, and on aggregate:

$$\sum_{t=0}^{\infty} q_t \Pi_t^D = \sum_{t=0}^{\infty} q_t \left( p_t^{ED} (\zeta_t K_t^D) - r_t^D p^D K_t^D - p_t^{fuel} D_t^E \right)$$

subject to emissions constraint (9), and a constraint on the utilization rate:  $\zeta_t \leq 1$ .

#### 5.5 The fossil-fuel-extracting firms' problem

Following Sinn (2008), we assume that the fossil-fuel extracting sector is competitive. Each representative firm possesses a fixed and known stock of the resource of the size  $s_0$ . Extraction costs depend negatively on the stock  $s_t$  remaining in the ground. The fossil fuel extraction cost per unit is given by  $\gamma_1 \left( \frac{s_0}{s_t} \right)^{\gamma_2}$ , from which it follows that aggregate extraction costs per unit may be written as

$$G^D(S_t) = \gamma_1 \left( \frac{S_0}{S_t} \right)^{\gamma_2}, \quad (10)$$

where  $S_t$  is the total stock remaining at time  $t$ , and  $\gamma_1$  and  $\gamma_2$  are parameters. This equation implies that when a smaller stock is left, the extraction cost will be higher. In this paper, we focus on coal: this represents the greater share of potential future emissions, and also is the relevant fuel for sectors such as coal-fired power stations, for which there are significant potential irreversibilities in investments. We calibrate (10) to this sector: coal is relatively abundant, with relatively flat extraction costs. See details in Appendix B.

Each individual firm extracts an amount of resource  $d_t^E + d_t^g$  in each period, which as usual we scale up across the economy to obtain the standard aggregate depletion equation:

$$S_{t+1} = S_t - (D_t^E + D_t^g).$$

The small firms act as price-takers, and the market price for their fossil fuel is  $p_t^{fuel}$ . By choosing an optimal extraction amount at each point in time, each firm in this sector maximizes the present value of its profits, and so on aggregate:

$$\sum_{t=0}^{\infty} q_t \Pi_t^{DE} = \sum_{t=0}^{\infty} q_t [p_t^{fuel} - \tau_t^D - G^D(S_t)] (D_t^E + D_t^g)$$

where  $\tau_t^D$  is a tax on the production of fossil fuels.

## 5.6 The renewable energy firms' problem

The renewable sector is composed of small firms, who do not internalize the learning-by-doing externality (5) (see discussion earlier). That is, on aggregate across these firms, they take the stock of accumulated knowledge about using the renewable energy  $H_t$  as given, with a rental rate  $r_t^H$ . They receive a subsidy of  $\tau_t^H$  on their dollar-valued holdings of renewable energy capital  $H_t$  and sell their output to the aggregate-electricity-producing firms at price  $p_t^{EH}$ . The firms take all prices as given, so on aggregate they maximize:

$$\sum_{t=0}^{\infty} q_t \Pi_t^H = \sum_{t=0}^{\infty} q_t [p_t^{EH} - p_t^H (r_t^H - \tau_t^H)] H_t.$$

Note that in the “simple model” of Section 3 we did not model renewable energy firms explicitly, so in that model we wrote the subsidy as accruing to the householder, who also owns the capital.

## 5.7 Climate system, emissions and damages

The carbon dioxide emissions  $D_t$  have three sources: “general” output production  $D_t^g$ ; electricity production from using fossil fuel  $D_t^E$ ; and land use  $D_t^{\text{land}}$ .

$$D_t = D_t^E + D_t^g + D_t^{\text{land}} \quad (11)$$

Land-use emissions  $D_t^{\text{land}}$  are set exogenously as by Cai et al. (2016). We use the climate system of Cai et al. (2016), which adapts the climate system of DICE2013 (Nordhaus, 2014a) to an annual time step. As this component of our model has been described extensively in the previous literature, we omit explanation here, and simply denote the mapping from emissions to temperature by:

$$T_t = \mathcal{W}_t(D_0, \dots, D_{t-1}) \quad (12)$$

where  $T_t$  is global atmospheric temperature change over pre-industrial levels,  $D_s$  is fossil-fuel-related pollution at time  $s < t$  and the warming function  $\mathcal{W}_t$  relates these two variables.

Finally, the damage factor for “DICE damages” is given by

$$\Omega(T_t) = \frac{1}{1 + \varsigma_1 T_t^{\varsigma_2}}, \quad (13)$$

where  $\varsigma_1$  and  $\varsigma_2$  are parameters. However, a great deal of discussion in real-world policy-making focuses on limiting global temperature changes to 2°C. We simulate this constraint by letting

$$\Omega(T_t) = \frac{1}{(1 + \varsigma_1 T_t^{\varsigma_2})(1 + \varsigma_3 (T_t/\varsigma_5)^{\varsigma_4})} \quad (14)$$

with a small positive parameter  $\varsigma_3 = 0.001$ , a large exponent parameter  $\varsigma_4 = 50$ , and  $\varsigma_5 = 2$ . Thus, when atmospheric temperature increase  $T_t$  is smaller than 2°C, the new damage factor given by (14) is almost the same as (13), but when  $T_t$  is larger than 2°C, the new damage factor will imply (much) larger damages than (13). This new damage factor (14) will be referred as the “stringent damage factor”.

Our primary motivation for using this specific stringent damage factor is that it ensures optimal policy in the model will indeed limit warming to 2°C. However, one should note that the damage factor (13), though commonly used, is very controversial (see, for example Weitzman 2009, 2010; Stern 2013; Pindyck 2013, 2017; Dietz and Stern 2015; Stern 2016; Cai and Lontzek Forthcoming).

In fact there do not exist well-founded estimates of global economic damages for even moderate temperature changes, and so the possibility to dictate optimal climate policy based on damage estimates is limited. Meanwhile, the view has been taken by many that warming above 2°C ought to be avoided. Such a constraint ought to be based on some premise as to the consequences of passing this threshold, which our damage factor provides.<sup>17</sup> Indeed, if recent warnings of the world entering a ‘hothouse Earth’ trajectory beyond 2 degrees of warming are to be believed, our stringent damage factor may not be unreasonable (see Steffen et al. 2018).

## 5.8 Decentralized equilibrium versus the social planner’s optimal solution

To find an optimal solution of the decentralized model, we formulate it as that of a principal who must choose an allocation from among those that can be implemented as a decentralized equilibrium, bearing in mind how the other economic participants (the “agents”) will respond. In the optimal taxation literature, such conditions imposed on the (Ramsey) principal are known as implementability conditions. We solve it using mathematical programming with equilibrium constraints. The details are in Appendix D.

The previous sections laid out the decentralized equilibrium model. To retrieve the values of the optimal carbon tax and optimal subsidies, that could replicate the first-best allocation in the decentralized equilibrium model, we also outline a social planner model where the social planner maximizes social welfare given constraints describing the carbon cycle, temperature, damages and fossil fuel depletion, and the capital accumulation equations. See Appendix C for details.

## 5.9 Subsidy and carbon tax

In the decentralized equilibrium, there are two instruments: a subsidy on renewable capital,  $\tau_t^H$ , and a carbon tax,  $\tau_t^D$ . There are various scenarios related to the choice of policy instruments. We differentiate between four cases: (1) a no policy scenario in which we set  $\tau_t^D = 0$  and  $\tau_t^H = 0$ ; (2) the optimal policy version, in which both instruments are freely chosen to maximize the principal’s objective; (3)  $\tau_t^D = 0$  and the subsidy is chosen freely to maximize the principal’s objective; (4)  $\tau_t^H = 0$  and the carbon tax is chosen freely to maximize the principal’s objective. Clearly, the second policy yields the same outcome as the social planner’s problem, and it is the first-best, which we prove in the appendices. Cases (3) and (4) are situations with second-best policies.

As is standard, we define:

**Definition 5.1.** The *social cost of carbon* (SCC),  $\chi_t$ , is the shadow price on carbon emissions, relative to the shadow value of output. That is, if  $\mu_t^D$  is the shadow price of Equation (11) constraining total emissions, then:

$$\chi_t := \frac{\mu_t^D}{u'(C_t/L_t)}.$$

Again we will write  $g_t^H = \frac{H_{t+1}-H_t}{H_t}$  and we prove (see Appendix D.8):

---

<sup>17</sup>One could alternatively use an explicit constraint that  $T_t \leq 2$  for all  $t$ . We prefer not to, in order to maintain continuity and smoothness in this part of the model. The explicit constraint  $T_t \leq 2$  creates difficulties in numerically solving for the decentralized equilibrium under the principal-agent framework (discussed in Section 5.8 and Appendix D) due to its associated complementarity condition (although it has no problem in solving the social planner’s problem). This will be compounded by the “irreversibility constraints” for the investments, especially (6). Moreover, in any scenario incorporating policies which can be used to limit temperature change at 2°C, an optimal economy facing our stringent damage factor (14) will very closely approximate the solutions we would obtain by using such a direct constraint. See also Appendix D.9.

**Proposition 5.2.** *The decentralized equilibrium allocation coincides with the solution to the social planner’s problem if carbon taxes are set as the social cost of carbon  $\chi_t$ , which is equal to:*

$$\chi_t = -u' \left( \frac{C_t}{L_t} \right)^{-1} \sum_{m=1}^{\infty} \beta^m u' \left( \frac{C_{t+m}}{L_{t+m}} \right) \frac{\partial Y_{t+m}}{\partial D_t}; \quad (15)$$

*and if subsidies are set equal to the “learning effect”:*

$$\tau_t^H = -(H_{t+1} - (1 - \delta^H)H_t) \frac{G'(H_t)}{p_t^H} = \lambda(g_t^H + \delta^H). \quad (16)$$

Equation (15) says the social cost of carbon is equal to be the marginal effect of present emissions on future welfare. This welfare impact is of course determined by the damage factor in use; we distinguish scenarios by the use of (13) and (14). Of course, our stringent damage factor (14) is designed to constrain emissions to 2°C, and so not to accurately depict true welfare damages. However, as discussed in Section 5.7 above, climate change economics has so far failed to identify a ‘true’ damage factor. Appendix D.8 provides further discussion of an alternative interpretation of (15) in the case of the stringent damage factor (14), as incorporating a shadow price of constraining temperatures to 2°C warming.

Equation (16) verifies that the theoretical insights on learning-by-doing from the “simple model” in Section 4 all carry across to the full model. That is, Corollary 4.3 holds and we have an “acceleration effect”. In particular, an increase in the carbon tax which reduces investment in and utilization of dirty energy capital and so increases deployment of the substitute renewable energy capital, *also* implies an increase in the optimal renewable energy subsidy.

Naturally, Proposition 5.2 also shows that we can examine optimal policy by using a social planner’s model, which is easier computationally. However, we do not restrict attention to this simpler case; we are also very interested in worlds without optimal (first-best) policy. If only the tax, or only the subsidy, are in use, then Proposition 5.2 does not apply. We explore such scenarios with our numerical results.

## 6 Quantitative results from the calibrated model

This section presents the quantitative results in four parts.<sup>18</sup> The first investigates the links between irreversible investment decisions and climate policies. We compare optimal policies with and without a stringent climate target (using the stringent damage factor (14) or the DICE damage factor (13) respectively). In addition, we illustrate the importance of the irreversibility in investment decisions relative to the case in which investments are reversible (the irreversibility effect). In the second part we study the acceleration effect pertaining to an early start of investment in the renewable sector. Next, we study the impact of climate policy stringency on the optimal carbon tax as well as the effect of learning-by-doing on the optimal carbon tax. Finally, we study the implications of the second-best policies for social welfare and the dynamics of the model.

The initial period in our model is 2012. Scenarios can be differentiated along three different dimensions: (1) damage function (DICE damage factor (13) vs. stringent damage factor (14)); (2) irreversible vs. reversible investments; and (3) the choice of policy instruments (optimal tax and subsidy vs. second-best policies). The runs of the decentralized equilibrium under the combined optimal tax and subsidy are equivalent to the runs of the social planner model (the first-best policy).

---

<sup>18</sup>The calibration of the model is described in Appendix B.

## 6.1 Irreversible investment and its implications

First, we want to understand how the optimal paths of variables depend on the irreversibility assumption coupled with different climate policy targets. We notice that the effect of irreversibility (compared with when the investment is reversible) becomes quantitatively important only if the climate policy objective is ambitious enough. Figure 1 shows that the paths of investment on dirty energy are almost the same with reversible and irreversible investments under a mild climate policy objective (the DICE damage factor) in this century, but they are distinctly different from each other under a more ambitious climate policy target (the stringent damage factor).

These results emphasize the importance of setting ambitious climate policies to induce permanent fuel energy switching. The strong path dependence embodied in carbon-intensive infrastructure suggests that mild climate change policies (those based on the DICE damage factor) would not induce fast shifts away from dirty energy towards green energy, as would be required to meet the Paris Agreement objectives.<sup>19</sup> (For the temperature pathways, see Figure 8 in Section 6.4.)

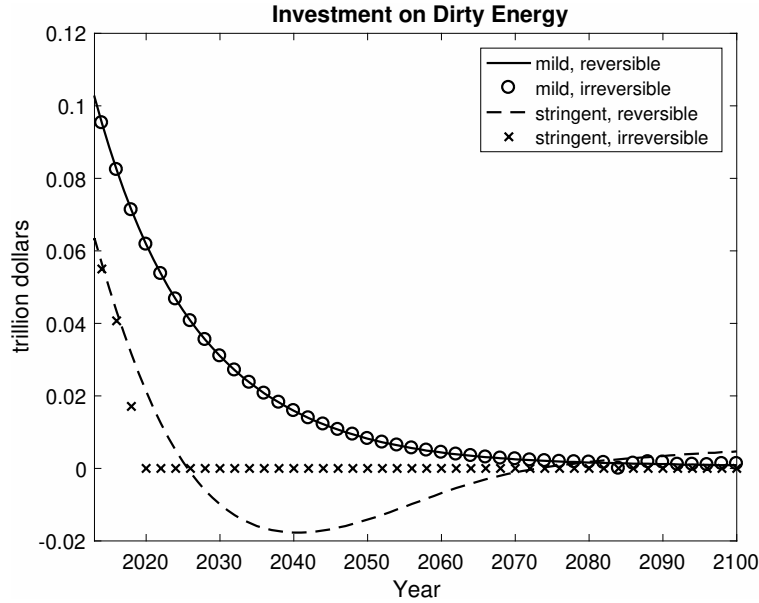


Figure 1: Investment in the Dirty Energy Capital Stock

Further, Figure 1 shows that with irreversible investments and the stringent policy target, there is no investment in dirty energy after 2020. In contrast, when investment is reversible, we keep investing in this capital stock until 2027 (another seven years). We then start turning dirty capital stock into general capital, a process that continues until about 2075. However, we never entirely stop using the dirty capital stock because of the imperfect substitutability between dirty and clean energy in electricity production (equation (8)). So, since we decumulated the dirty capital stock sufficiently in the preceding decades, investment in the dirty capital stock resumes after 2075 under reversible investment.

These dynamic patterns of investment in dirty energy with the (ir)reversible investments and the stringent damage factor correspond to the dynamics of returns on those investments shown in Figure

<sup>19</sup>This finding echoes the one in Meng (2016), who estimates the strength of path dependence in the electricity sector for the U.S. Midwest and shows that a permanent decline in U.S. electricity sector emissions would require shocks of larger magnitude and longer duration than that of recent natural gas prices.

2.<sup>20</sup> The theoretical counterpart of this figure is Lemma 3.2 in Section 3. First, the figure shows that we end investment in the dirty capital stock when the investment is still attractive with the rate of return,  $r_t^D - \delta^D$ , exceeding the rate of return on the general economy,  $r_t^g - \delta^g$ . This is because we will only invest in infrastructure that will become obsolete if the short-term benefits from that investment compensate for future losses. Thus, even without uncertainty, returns to irreversible investment require a premium.<sup>21</sup> Even if we end investment at around 2020, we continue to fully utilize the dirty capital stock for about another 25 years, until 2045, when the return on dirty capital ( $r_t^D$ ) reaches zero and we start underutilizing the dirty capital stock.

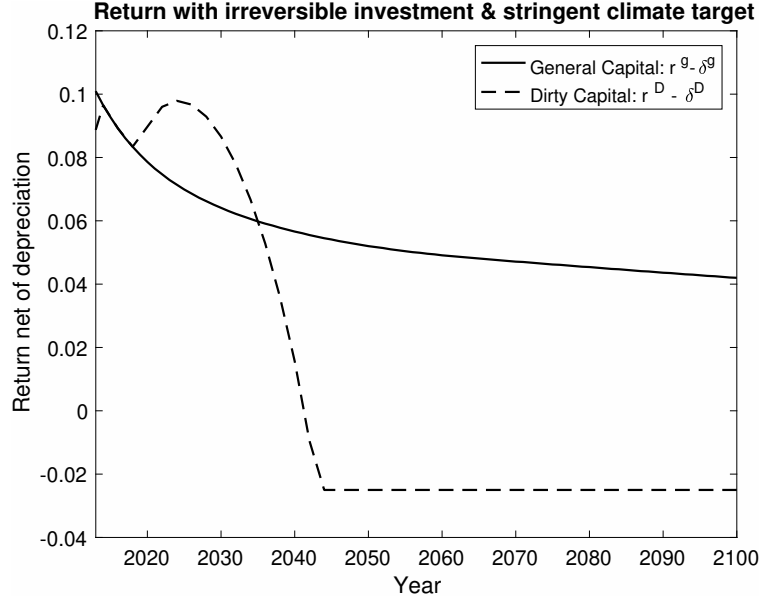


Figure 2: Return on general and dirty capital

Consider the medium-term, until about 2037. As we end investment sooner than in the counterfactual (when disinvestment is a viable option), the economy continues to hold a smaller amount of dirty capital stock under irreversibility compared to the reversible case (the solid and dashed lines in Figure 3a). After 2037, however, the economy holds larger stocks of dirty capital in the long-run if investment is irreversible. This result is due to path dependence: capital cannot be converted into other forms of capital stock. However, if we take into consideration the underutilization of the dirty capital stock in the irreversible investment case (marks in Figure 3a), then in the long-run, the same total amount of dirty capital stock will be *utilized* under irreversible and reversible investment decisions (marks and circles in Figure 3a).

Because emissions from the dirty energy sector are directly proportional to utilized dirty capital,

<sup>20</sup>The small divergence between returns in the first period is a standard artifact in dynamic models with given initial conditions: our calibration provides a little more dirty capital stock than the short-run equilibrium would dictate, and so returns are below the general economy for one period. A pause in investment brings the model into balance.

<sup>21</sup>This finding has empirical support from previous studies on irreversible investment in other contexts. For example, Bernstein and Mamuneas (2007) develop a simple model of production and investment with costly disinvestment to estimate the magnitude of the premium associated with irreversible investment in the telecommunications industry, assuming future telecommunications capital acquisition prices are random variables. Their findings indicate that the premium increases the user cost of capital by 70%, which implies an average hurdle rate of 14% over the period 1986-2002. Using different methods and framework, Pindyck (2005) provides similar estimates of the telecommunications hurdle rate.

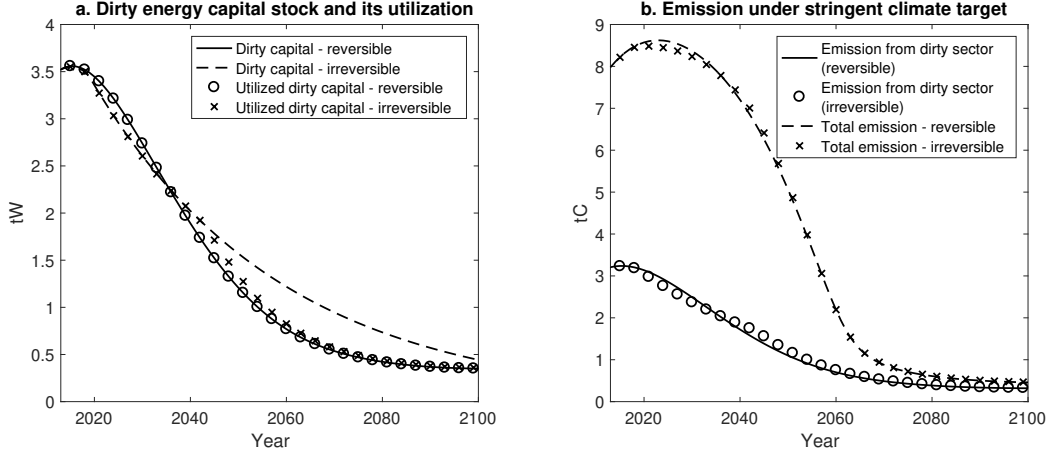


Figure 3: Dirty capital and emissions under the stringent climate policy target.

the utilization curves in Figure 3a also give the pattern of emission levels. We show the level of these emissions, as well as total emissions (i.e.,  $D_t$ , including those from the general economy and land-use) in Figure 3b.

## 6.2 The Acceleration Effect

The theoretical result in Section 4.2 gave the optimal subsidy, when the optimal carbon tax is also present and when the irreversibility constraint on investment in this sector does not bind (Corollary 4.3):  $\tau_t^H = \lambda (g_t^H + \delta^H)$ . This formula implies that (i) the subsidy continues as long as there is investment in the renewable sector, (ii) the subsidy increases with the learning coefficient  $\lambda$ , and (iii) the optimal subsidy is higher when optimal growth in renewable capital is higher. There are four different mechanisms that can lead to higher capital accumulation in the renewable sector and consequently to a higher level of subsidies: (1) more stringent climate policy targets; (2) the dirty sector could be shrinking faster than in the reversible case due to the irreversibility effect; (3) a higher learning rate (from higher R&D in renewables); and (4) under second-best scenario when a carbon tax is not possible, it could be optimal to grow the renewable sector faster to crowd out the dirty energy sector. We here investigate the first three of these channels, noting that the irreversibility constraint on investment in renewables never binds in our model. We will consider second-best policies, which encompass many important effects, in Section 6.4.

### 6.2.1 Channel 1: stringent climate policy

Figure 4a plots the optimal subsidy  $\tau_t^H$ , and Figure 4b plots the total subsidy (i.e., total amount of dollars paid),  $p_t^H \tau_t^H H_t$ , under mild and stringent climate policy targets. In the stringent case we plot with both reversible and irreversible investments on the dirty side of the energy economy, as the clean and dirty sides will interact.

We observe in Figure 4a that the subsidy is higher under the stringent climate policy target in the initial decades (following 2012). Because we use less fossil fuels in this scenario, we must generate more of our electricity from renewables, and so the latter sector is initially growing faster than it is in the mild policy scenario. From mid-century onward, this order reverses: if we ignore the additional effects due to irreversibility, a higher subsidy is paid under the mild policy target. This is because, in the stringent policy scenario, a substantial volume of renewable energy capital

stock has already been built by this point, and thereafter its optimal growth rate is slower; by the acceleration effect, so also is the optimal subsidy.

We also observe these effects when we plot the *total* subsidy paid to all holders of renewable energy capital  $H_t$ , in Figure 4b. This shows that payments are always higher under the stringent policy target, even when the subsidy (and indeed price) of the capital stock are lower, because of the size of  $H_t$ . (The decline in growth in these subsidy payments, starting at around 2040, mirrors the decline in growth in  $H_t$  already seen at this time). Still, the total subsidies of the three scenarios are converging as we approach the end of the century.

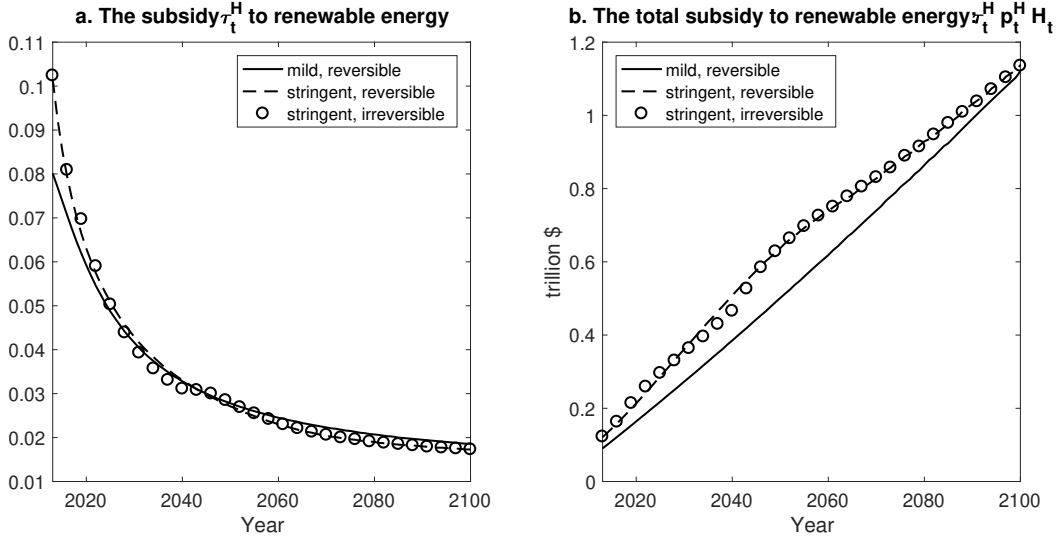


Figure 4: Optimal Subsidies under Our Three Main Scenarios.

### 6.2.2 Channel 2: interaction with the irreversibility effect

In Figure 4, we plotted the optimal subsidies under the stringent climate policy with both *irreversible* and *reversible* investment decisions. There is a complex relationship between the subsidy level and the results of the irreversibility effect discussed in Section 6.1. As seen in Figures 1 and 3a, in the early periods the dirty sector shrinks faster in the irreversible case, as compared with the reversible case. However, this pattern reverses so that after around 2035 there is more dirty energy capital in use in the irreversible case. Underutilization begins around 2040, but there is still excess dirty energy capital in use until around 2070, when the trajectories converge.

So in those earlier years, when we build fewer coal-based power plants, we must be building a greater volume of renewables instead. This explains the greater subsidy for renewables in the irreversible case, visible in Figure 4a up to around 2020.

However, after this point, the subsidy to renewables drops below that for the reversible case. This is in anticipation of the greater dirty energy capacity that will remain in the economy, due to irreversibility (instead of being decumulated, as in the reversible case). Thus, until 2045, the subsidy in the irreversible case remains below that of the reversible case—and even below that of the mild policy target.

But as the use of dirty capital begins to approach that of the irreversible case, it becomes necessary to again accelerate the deployment of the substitute renewable capital. Thus, from 2045, the subsidy to renewables in the irreversible case again exceeds that of the reversible case. The



total subsidies paid are approximately the same in both cases at this point (Figure 4b) because less of this capital stock has been accumulated in the irreversible case, due to the prolonged reduction in investment.<sup>22</sup>

### 6.2.3 Channel 3: learning rate

Figure 5 presents the level of subsidies for three different years (2018, 2050 and 2100) under different values of the learning parameter  $\lambda$ . There are two important features to note. First, the level of subsidies increases with the value of learning parameter. Second, the subsidy level for 2018 follows a convex pattern with respect to the value of  $\lambda$ , whereas it follows a more linear pattern for the other two years.

Both of these patterns can be explained by referring to our theoretical result:  $\tau_t^H = \lambda (g_t^H + \delta^H)$  (Corollary 4.3). The primary increase of subsidy with  $\lambda$  is clear. In addition, the degree of convexity of the optimum subsidy level will be determined by the optimal growth of deployment of renewables,  $g_t^H$  (our acceleration effect). Specifically, for 2018, there is higher optimal growth of renewables with a higher learning rate, because this higher rate makes it more economically advantageous to expand the sector. Thus, since a higher value of  $\lambda$  means that  $g_{2018}^H$  is *also* higher, it follows that  $\tau_{2018}^H$  is convex with  $\lambda$ . However, in 2050 and 2100, the higher growth in renewables has already taken place and the renewables are already functioning in the economy as mature technologies. As such there is no need to grow the renewable sector as fast in 2050 and 2100 as in 2018. Moreover, a higher learning rate, and so faster growth in the earlier periods, lead to greater market saturation later in the century. Thus, at this time, optimal growth is slightly lower in the case of higher  $\lambda$ . Thus, the relationship between  $\lambda$  and the becomes slightly concave at 2050 and 2100.

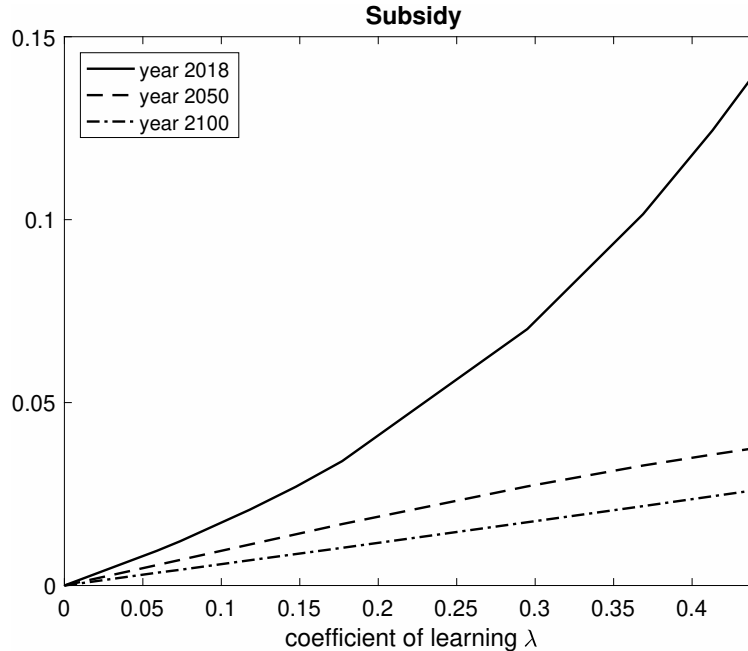


Figure 5: Optimal subsidy for different values of learning rate  $\lambda$  in the reversible and stringent scenario.

<sup>22</sup>Related literature has investigated the optimal time path for innovation policy, see, for example Gerlagh et al. (2009) and Gerlagh et al. (2014). For instance, the latter show that if the patent lifetime is finite, the optimal subsidy starts at a high level, providing an incentive to accelerate R&D investments, and then falls over time.

### 6.3 Optimal Carbon Taxes

The social cost of carbon (SCC) is generally considered to be the most important concept in climate change economics. If all other externalities are internalized, the optimal carbon tax should be set to this level (Proposition 5.2). It is equal to the marginal effect on future welfare of present emissions, via their effect on output, normalized according to the marginal utility of individual consumption at the time at which it applies.

Figure 6 displays the impact of climate policy stringency on this tax, for the cases of reversible investments, where stringency is measured by the parameter  $\varsigma_5$  in (14). We use three levels of stringency:  $\varsigma_5 = 2$  is the default stringent climate policy (the “stringent” policy target used in previous sections),  $\varsigma_5 = 2.2$  the medium stringency, and  $\varsigma_5 = 2.5$  the low stringency. We also include results for the mild policy target (as defined in previous sections) which is equivalent to the case with  $\varsigma_5 = \infty$ . The impact of stringency on carbon tax is large. For example, in 2050, the optimal carbon tax under the mild, low stringency, medium stringency, or the default stringent climate policy is \$226/tC (USD per ton of carbon), \$261/tC, \$386/tC, and \$542/tC, respectively.

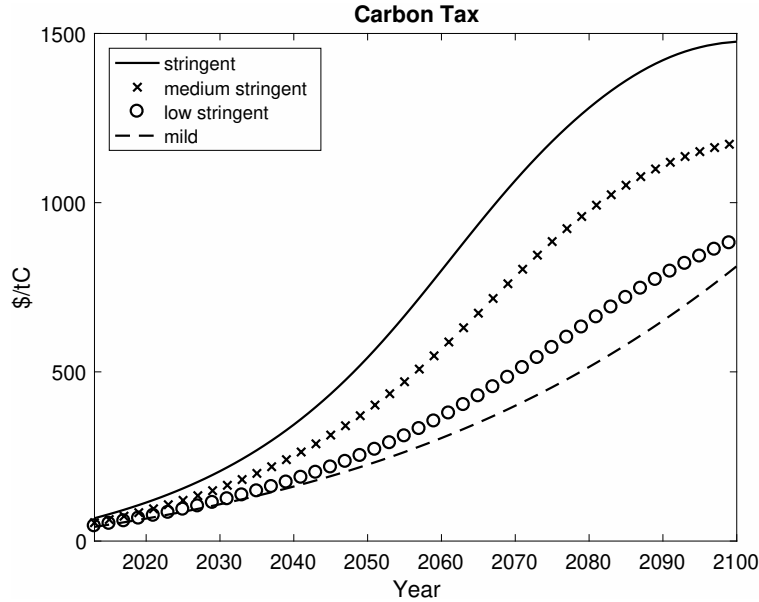


Figure 6: The Effect of Policy Stringency on the Optimal Carbon Tax

Figure 7 presents the effect of learning-by-doing on the carbon tax in the case with reversible investments: with a stringent climate policy target, the carbon tax is higher *without* learning-by-doing as it could be expected. With learning-by-doing, the associated subsidized roll-out of renewable energy technologies means that emissions are lower, both in the current period and in the future. It follows, in the stringent scenario, that the marginal effect on future welfare of current emissions is lower: a lower carbon tax is optimal. Intuitively, this is because cheap low-carbon energy means that stringent policy targets can be met without imposing a higher carbon tax.

### 6.4 Second-best policies

#### 6.4.1 Subsidy Versus Tax

As Proposition 5.2 showed, the decentralized equilibrium with the optimal carbon tax on the externality created by fossil fuel use in the energy sector, combined with the optimal subsidy on the

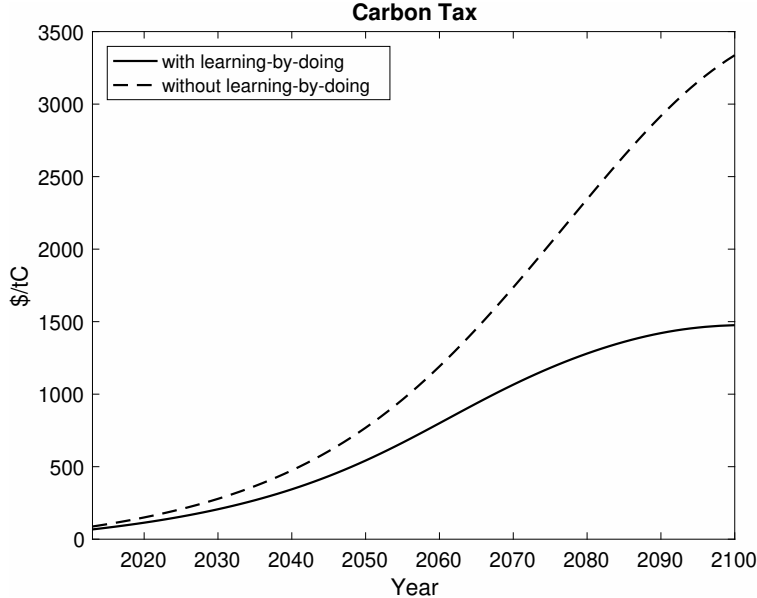


Figure 7: Optimal Carbon Tax with/without learning-by-doing in the stringent and reversible scenario.

learning-by-doing externality in the renewable sector, implements the optimal allocation obtained in the social planner’s problem (the first best). In practice, however, one of these two policy instruments may be unavailable, and policy makers thus have to rely on second-best policies. In this section we compare the relative performance of these two policy instruments when used alone, under alternative climate policy objectives and (ir)reversible investment decisions. This is an important exercise given the landscape of debate regarding optimal climate policy. While these second-best policies represent two extremes, considering these extremes gives us the extent of the differences, and results for intermediate policies may be interpolated (we consider an intermediate case in Section 6.4.3 below).

Moreover, some criticize existing subsidies as expensive and inefficient, and would prefer something closer to a tax-only policy (see e.g. Helm 2012). Conversely, if we look at actual implementations, we see that the European Union has spent a large amount on subsidies to achieve its target of 20% renewable energy by 2020, while allowing the carbon price in its trading scheme to fall to very low levels. In between, many advocate the necessity of mixed policies, while stressing the critical importance of carbon pricing.<sup>23</sup>

We contribute to this debate, arguing that in a second-best world, the policy instrument that should be used depends on how stringent climate policy objectives are. More specifically, under mild climate policy targets, as in the case with the ‘DICE’ damage factor (13), the economy is better off (in the sense of social welfare) with the optimal subsidy as a policy instrument. In contrast, under more stringent climate policy targets, as in case with the stringent damage factor (14), the economy is better off if the optimal carbon pricing policy is adopted (see Table 1).<sup>24</sup> As the results reported in the Table 1 further indicate, irreversibility in investment decisions does not affect the

<sup>23</sup>Bowen (2011) argues that “other policies are needed, too, particularly to promote innovation and appropriate infrastructure investment, but cannot be relied upon by themselves to bring about the necessary reductions to emissions. Carbon pricing is crucial.”

<sup>24</sup>These findings are in line with ones in Gerlagh and van der Zwaan (2006) who use a long-term top-down model, with a decarbonization option through carbon capture and storage, to show that carbon taxes do better for stringent targets, and subsidies do better for modest targets. However, we use a different approach from them, analyzing the

|  | Optimal tax<br>zero subsidy | Optimal subsidy<br>zero tax |
|--|-----------------------------|-----------------------------|
| Reversible investment<br>mild climate policy target        | 2.05%                       | 1.29%                       |
| Irreversible investment<br>mild climate policy target      | 2.05%                       | 1.19%                       |
| Reversible investment<br>stringent climate policy target   | 2.55%                       | 3.96%                       |
| Irreversible investment<br>stringent climate policy target | 2.55%                       | 2.94%                       |

Table 1: Second-best policies: welfare loss, % of initial period consumption

relative ranking of these policy instruments. Interestingly, incorporating irreversibilities *reduces* the percentage loss from using the subsidy, under both policy scenarios, while making little difference to the loss from the tax. We now discuss the reasons for these results.

Figure 8 shows the temperature, emission, and tax levels under mild climate policy targets (the left panels) and stringent climate policy targets (the right panels), both assuming irreversible investments. The top-left and middle-left panels show that with only carbon pricing, temperature and emissions paths closely follow those under the first-best policy. This is accomplished with a (slightly) higher level of carbon tax than under the first-best scenario. If we consider the more stringent climate policy case, we observe a similar pattern of paths for temperature and emissions: with carbon pricing only, they closely follow the paths of the first-best (the top-right and middle-right panels of Figure 8). The second-best tax level is again higher than the first-best counterpart.

The intuition behind these results is as follows. With only carbon pricing, there is a risk of being “locked into” the ways of producing electricity that are currently cheap: coal-based power plants.<sup>25</sup> Meanwhile, the alternative (producing electricity from renewables) is currently more expensive and may not become competitive. As a result, a higher level of carbon taxes on the fossil fuel extracting firms is needed compared with the first-best scenario. But since the size of the dirty sector in the energy sector of the economy is large relative to the renewable sector, this policy of making the dirty sector “less competitive” through carbon taxes is relatively more costly (in welfare terms), than the policy of making the renewable sector competitive through direct subsidies.

But, unlike carbon pricing, subsidies directly stimulate investment in renewable energy and, once clean technologies develop and become competitive, the renewable sector crowds out the dirty energy sector. Under less ambitious climate policy, this subsidy appears sufficient, as well as less costly than carbon pricing (given also the relatively smaller size of the clean sector). On the other hand, achieving the more stringent climate policy target through innovation policy is extremely difficult as it requires decarbonization of the large dirty energy sector. Adoption of the instrument which directly targets that sector—carbon pricing—is a policy that is associated with relatively higher welfare.

Finally, introduction of irreversibilities reduces welfare both in the first- and second-best cases, as the irreversibility constraint (6) is always binding at some periods (in the mild policy case it starts to bind around the year 2080, and so this is not clearly visible in Figure 1). But, because the emission paths, and investments in dirty capital, are very similar in the “tax-only” and “optimal” scenarios, the introduction of these irreversibilities makes no discernible difference to the welfare

implications of the second-best instruments for climate policy in a transparent setting.

<sup>25</sup>See Unruh (2002) and Jaffe et al. (2005) for further discussion.

lost relative to the first best, when we use the tax-only policy. This is shown in Table 1.

However, irreversibilities change the optimal subsidy: our acceleration effect (Sections 4 and 6.2.2). So, if the only policy available is a subsidy, then re-optimization by the policy-maker means that the subsidy and renewable investment levels are higher and earlier, in the presence of irreversibilities. This change in the investment trajectory, capturing the benefits of early learning-by-doing, reduces the net impact of the irreversibility constraint. Thus, although welfare has been lost due to irreversibilities, the *relative* losses are smaller than they were before, in the subsidy-only case: see Table 1.

Finally, the emissions and temperature paths with carbon pricing only, irrespective of the assumptions about the stringency of climate policy targets, closely follow the ones of the first-best scenario because carbon pricing internalizes the global warming externality, and thus is better suited to target climate policy objectives.

#### 6.4.2 Subsidy versus tax: discussion

Our results on the relative performance of carbon pricing versus subsidies in a second-best setting reflect the broad trends in the global climate political landscape. Nowadays we observe a rapid expansion in the use of renewable energy technologies.<sup>26</sup> Renewable energy technologies are viewed today as tools to mitigate climate change, improve local air quality, advance economic development and create jobs. Declining costs have played a pivotal role in the expansion of renewable energy technologies in recent years. The stage for such an expansion was set more than a decade ago when a handful of countries, including Germany, Denmark, Spain, and the United States, created a critical market for renewables, which drove early economies of scale and led to the changes we witness today (Lins et al., 2014). During that period and effectively until 2016, when the Paris Agreement came into force, progress in the area of international climate policy had been modest at best. Although the European Union had started campaigning for the 2°C target in the mid-1990s, this target was not formally adopted until 2010 at the UN Climate Change Conference in Cancun. As such, we could characterize the international climate policy up to 2015 as having unambitious climate policy objectives. The Paris Agreement, however, renewed the climate political landscape, at least in theory, with a larger recognition of the urgency of more ambitious emissions reductions. The agreement has also revived discussion about the importance of adopting carbon pricing to implement the emissions mitigation pledges submitted by 186 countries for the December 2015 Paris Agreement,<sup>27</sup> which is in line with the message from simulations of our model under the second-best setting that more ambitious climate policy should adopt carbon pricing.

#### 6.4.3 Optimal Subsidy with a Pre-specified Tax

Section 6.4.1 analyzed two extreme cases. In practice, an intermediate situation may hold. It may be possible to have both a subsidy, and a carbon tax—but it may be politically impossible to set the tax as high as its optimal level. Policy-makers must then adjust the subsidy to meet the policy target. For example, the carbon tax could be pre-specified as half of the optimal carbon tax, where the optimal carbon tax is the solution of the first-best policy with both the carbon tax and subsidy are available. We use the mild target and reversible investment scenario to study this intermediate situation. Figure 9b displays the pre-specified carbon tax, and Figure 9a provides the corresponding

<sup>26</sup>Renewables accounted for nearly half of all new power generation capacity in 2014, led by growth in China, the United States, Japan and Germany, with costs continuing to fall (EIA, 2015).

<sup>27</sup>Baranzini et al. (2017) provide a summary of the main arguments in favor of carbon pricing in a post-Paris world. See also Farid et al. (2016) who urge for carbon taxes (or equivalently carbon trading systems) for implementation of the Paris pledges.

optimal subsidy. We see that when the tax is set to be zero, the subsidy is the largest; when the tax is set at its optimal level under the first-best policy, the subsidy is the smallest; and when the tax is set to be half of the first-best optimal level, the subsidy is between the previous two extreme cases. However, this pattern is only valid before 2055: after 2070, a smaller tax will be accompanied by a smaller subsidy, as larger subsidies in the earlier periods lead to a higher renewable energy stock, lessening the need for subsidies in the later periods.

## 7 Conclusion

In this paper we have studied implications of two capital stock effects – path dependence in infrastructure and learning-by-doing in the renewable sector – for the design of optimal climate policies, using both simple analytical models and simulations of the full dynamic general equilibrium climate-economy model.

We find that for temperature changes to not exceed 2°C, investments in dirty infrastructure should end in 2020. We show that the “Green Paradox” – that future stringent climate policy raises short-term emissions – has a converse if we focus on demand side capital stock effects. If the dirty capital stock cannot be converted to other forms of capital, then it is optimal to stop investing in the dirty capital stock earlier than the case where capital investments are reversible. Learning-by-doing significantly advances the timing of investment in renewables, not only to prevent later stranding of fossil-fuel-based assets but also to accelerate the decline in the costs of clean energy.

The timing of these effects depends on the stringency of climate policy targets. Climate policy targets induce an earlier shift (within the next few decades), to clean energy and away from dirty energy, only if they are stringent. Otherwise, path dependence in energy systems and low substitutability between the dirty and clean energy sources imply a prolonged period of using the dirty capital stock.

We have investigated carbon taxes and renewable subsidies as policy instruments in a perfect foresight, deterministic setting. Incorporating uncertainty, or wider issues with implementation costs and political economy, would invite comparisons with alternative quantity-based policies, such as “cap and trade” or “renewables obligations”. And indeed, risks such as climate tipping points, and uncertainties in damage functions, and advances of renewable energy technology, as well as carbon capture and storage, may have significant impact. These could be an interesting extension of the paper, which we leave for future research.

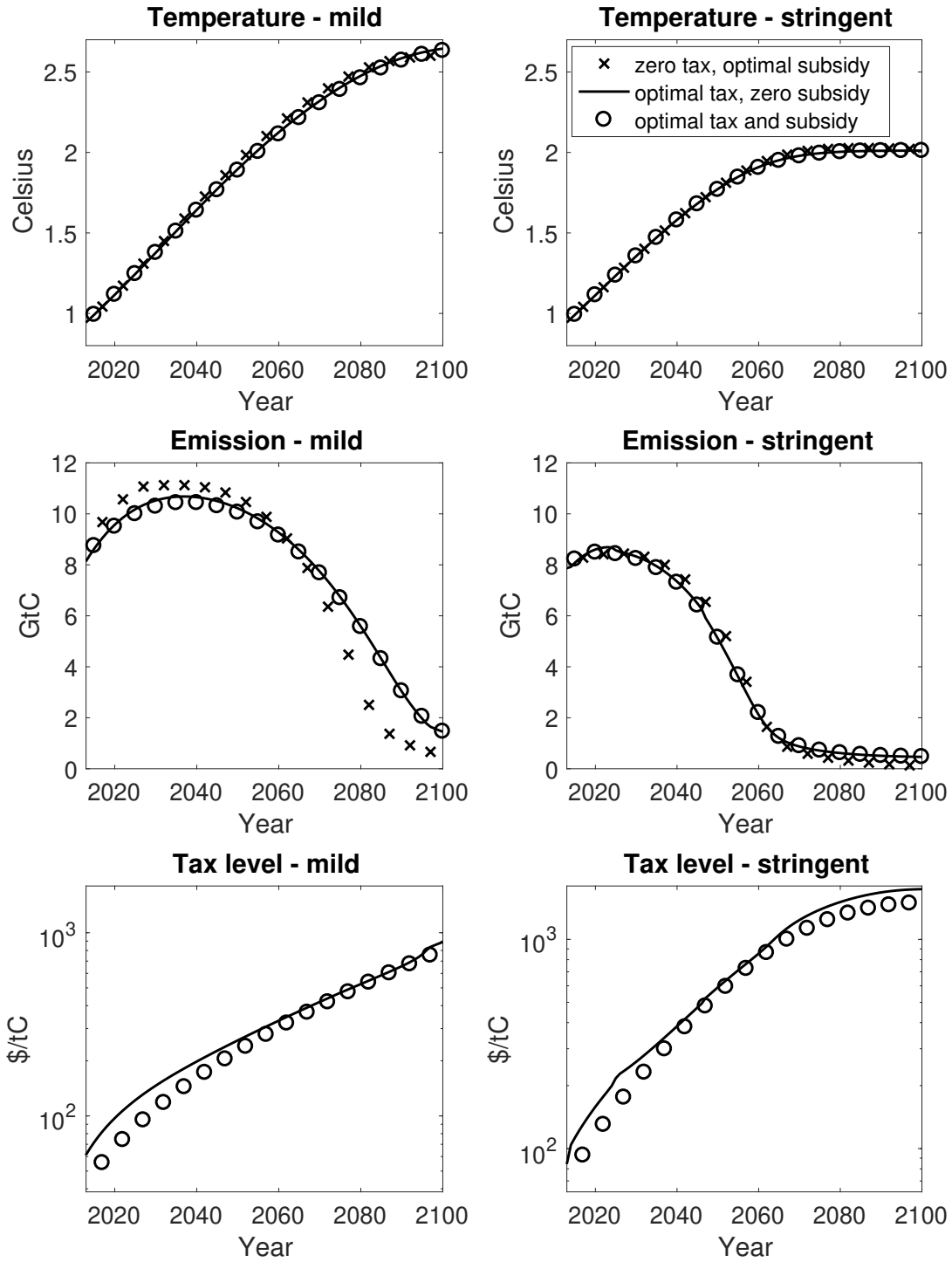


Figure 8: Temperature, Emission, and Tax level under the mild or stringent climate policy targets

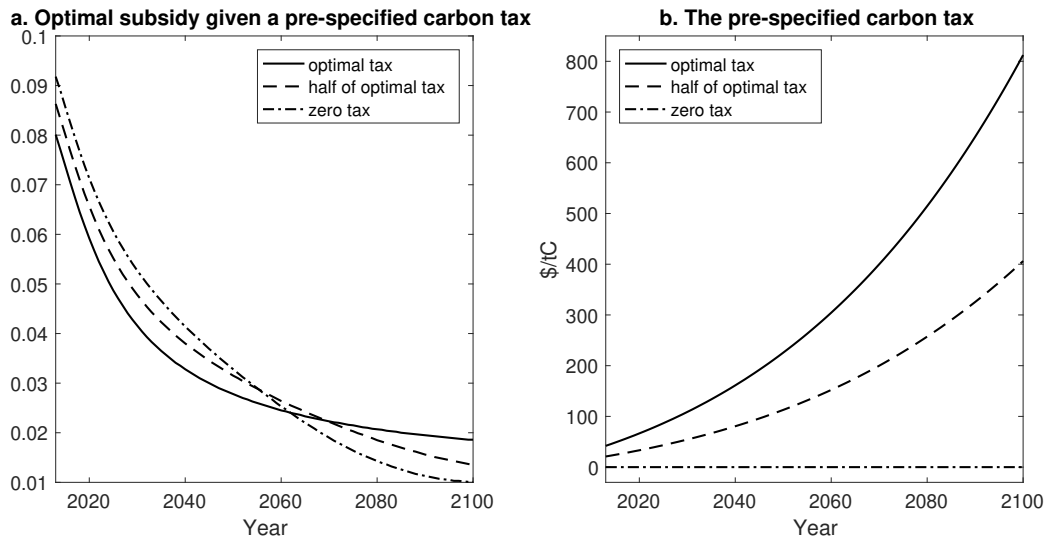


Figure 9: Optimal subsidy when the level of tax is constrained.



## References

- A. B. Abel. Optimal investment under uncertainty. *American Economic Review*, 73(1):228–233, 1983.
- D. Acemoglu, U. Akcigit, D. Hanley, and W. Kerr. Transition to clean technology. *Journal of Political Economy*, 124(1):52–104, 2016.
- P. Aghion, C. Hepburn, A. Teytelboym, and D. Zenghelis. Path dependence, innovation and the economics of climate change. *New Climate Economy contributing paper*, 2014.
- P. Aghion, A. Dechezleprêtre, D. Hémous, R. Martin, and J. V. Reenen. Carbon taxes, path dependency, and directed technical change: Evidence from the auto industry. *Journal of Political Economy*, 124(1):1–51, 2016.
- K. Arrow. The economic implications of learning by doing. *The Review of Economic Studies*, 29: 155–173, 1962.
- K. Arrow. Optimal capital policy with irreversible investment. In J. N. Wolfe, editor, *Value, Capital and Growth, Essays in Honor of Sir John Hicks*. Edinburgh: Edinburgh University Press, 1968.
- K. Arrow and M. Kurz. Optimal growth with irreversible investment in a Ramsey model. *Econometrica*, 38(2):331–344, 1970.
- A. Atkeson and P. Kehoe. Models of energy use: putty-putty versus putty-clay. *American Economic Review*, 89(4):1028–1043, 1999.
- A. Baranzini, J. C. J. M. van den Bergh, S. Carattini, R. Howarth, E. Padilla, and J. Roca. Carbon pricing in climate policy: seven reasons, complementary instruments, and political economy considerations. *WIREs Climate Change*, 8:1–17, 2017.
- L. Barrage. Optimal dynamic carbon taxes in a climate-economy model with distortionary fiscal policy. *Review of Economic Studies*, Forthcoming.
- N. Bauer, C. McGlade, J. Hilaire, and P. Ekins. Divestment prevails over the green paradox when anticipating strong future climate policies. *Nature Climate Change*, 8(2):130–134, 2018.
- J. I. Bernstein and T. P. Mamuneas. Irreversible investment, capital costs and productivity growth: Implications for telecommunications. *Review of Network Economics*, 6(3), 2007.
- A. Bhattacharya, J. Oppenheim, and N. Stern. Driving sustainable development through better infrastructure: Key elements of a transformation program. Working Paper 91, Global Economy and Development, July 2015.
- B. Bollinger and K. Gillingham. Learning-by-doing in solar photovoltaic installations. Available at SSRN: <https://ssrn.com/abstract=2342406> or <http://dx.doi.org/10.2139/ssrn.2342406>, 2014.
- A. Bowen. The case for carbon pricing. *Policy Brief, Grantham Research Institute on Climate and the Environment*, 2011.
- Y. Cai and T. S. Lontzek. The social cost of carbon with economic and climate risks. *Journal of Political Economy*, Forthcoming.

- Y. Cai, T. M. Lenton, and T. S. Lontzek. Risk of multiple interacting tipping points should encourage rapid CO<sub>2</sub> emission reduction. *Nature Climate Change*, 6(5):520–525, 2016.
- Y. Cai, W. Brock, A. Xepapadeas, and K. L. Judd. Climate policy under cooperation and competition between regions with spatial heat transport. NBER Working Paper 24473, 2018.
- B. Caldecott. Introduction to special issue: stranded assets and the environment. *Journal of Sustainable Finance & Investment*, 7:1:1–13, 2017.
- B. Caldecott, N. Howarth, and P. McSharry. Stranded assets in agriculture: Protecting value from environment-related risks. *Smith School of Enterprise and the Environment, University of Oxford*, 2013.
- G. Casey. Energy efficiency and directed technical change: implications for climate change mitigation. *Unpublished manuscript, Brown University*, 2017.
- S. Davis, K. Caldeira, and H. D. Matthews. Future CO<sub>2</sub> emissions and climate change from existing energy infrastructure. *Science*, 329:1330–1333, 2010.
- S. Dietz and N. Stern. Endogenous growth, convexity of damage and climate risk: How Nordhaus’ framework supports deep cuts in carbon emissions. *The Economic Journal*, 125(583):574–620, 2015.
- A. Dixit. Investment and hysteresis. *Journal of Economic Perspectives*, 6(1):107–132, 1992.
- EIA. Energy and climate change: World energy special report. *Report*, 2015.
- Energy and Environmental Economics, Inc. Generation cost model for China, December 2012.
- M. Farid, M. Keen, M. Papaioannou, I. Parry, C. Pattillo, and A. Ter-Martirosyan. After Paris: Fiscal, macroeconomic, and financial implications of climate change. *IMF Staff Discussion Note, SDN/16/01*, 2016.
- C. Fischer and R. Newell. Environmental and technology policies for climate mitigation. *Journal of Environmental Economics and Management*, 55(2):142–162, 2008.
- C. Fischer, L. Preonas, and R. Newell. Environmental and technology policy options in the electricity sector: are we deploying too many? *Journal of the Association of Environmental and Resource Economics*, 4(4):959–984, 2017.
- R. Fouquet. Trends in income and price elasticities of transport demand (1850-2010). *Energy Policy*, 50:62–71, 2012.
- R. Fouquet. Path dependence in energy systems and economic development. *Nature Energy*, N. 16098, 2016.
- D. Fudenberg and J. Tirole. Learning-by-doing and market performance. *The Bell Journal of Economics*, 14(2):522–530, 1983. ISSN 0361915X. URL <http://www.jstor.org/stable/3003653>.
- R. Gerlagh. Too much oil. *CESifo Economic Studies*, 57(1):79–102, 2011.
- R. Gerlagh and van der Zwaan. Options and instruments for a deep cut in CO<sub>2</sub> emissions: carbon dioxide capture or renewables, taxes or subsidies? *The Energy Journal*, 27(3):25–48, 2006.

- R. Gerlagh, S. Kverndokk, and K. E. Rosendahl. Optimal timing of climate change policy: interaction between carbon taxes and innovation externalities. *Environmental and Resource Economics*, 43(3):369–390, 2009.
- R. Gerlagh, S. Kverndokk, and K. E. Rosendahl. The optimal time path of clean energy R&D policy when patents have finite lifetime. *Journal of Environmental Economics and Management*, 67(1):71–88, 2014.
- M. Golosov, J. Hassler, P. Krusell, and A. Tsyvinski. Optimal taxes on fossil fuel in general equilibrium. *Econometrica*, 82(1):41–88, 2014.
- L. H. Goulder and K. Mathai. Optimal CO<sub>2</sub> abatement in the presence of induced technological change. *Journal of Environmental Economics and Management*, 39:1–38, 2000.
- J. Greenwood, Z. Hercowitz, and P. Krusell. Long-run implications of investment-specific technological change. *The American Economic Review*, 87(3):342–362, 1997.
- M. Grubb, T. Chapuis, and M. Ha-Duong. The economics of changing course: Implications of adaptability and inertia for optimal climate policy. *Energy Policy*, 23 (4-5):417–431, 1995.
- J. Hassler, P. Krusell, and C. Olovsson. Energy-saving technical change. Working Paper 18456, National Bureau of Economic Research, October 2012.
- D. Helm. *The Carbon Crunch: How We’re Getting Climate Change Wrong—and how to Fix it*. Yale University Press, 2012.
- A. B. Jaffe, R. G. Newell, and R. N. Stavins. A tale of two market failures: Technology and environmental policy. *Ecological Economics*, 54:164–174, 2005.
- S. Jensen, K. Mohlin, K. Pittel, and T. Sterner. An introduction to the green paradox: the unintended consequences of climate policies. *Review of Environmental Economics and Policy*, 9:246–265, 2015.
- D. Jorgenson. The theory of investment behavior. In R. Ferber, editor, *Determinants of Investment behavior*, pages 129–175. NBER, 1967.
- S. Kverndokk and K. E. Rosendahl. Climate policies and learning by doing: impacts and timing of technology subsidies. *Resource and Energy Economics*, 29:58–82, 2007.
- F. Lafond, A. G. Bailey, J. D. Bakker, D. Rebois, R. Zadourian, P. McSharry, and J. D. Farmer. How well do experience curves predict technological progress? A method for making distributional forecasts. *Technological Forecasting and Social Change*, 128:104–117, 2018.
- A. Lindman and P. Soderholm. Wind power learning rates: a conceptual review and meta-analysis. *Energy Economics*, 34(3):754–761, 2012.
- C. Lins, L. Williamson, S. Leitner, and S. Teske. The first decade: 2004–2014: 10 years of renewable energy progress. *Renewable Energy Policy Network for the 21st Century.*, 20:1–48, 2014.
- C. McGlade and P. Ekins. The geographical distribution of fossil fuels unused when limiting global warming to 2°C. *Nature*, 517:187–190, 2015.
- K. C. Meng. Estimating path dependence in energy transitions. Working Paper 22536, National Bureau of Economic Research, August 2016.

- T. Michielsen. Brown backstops versus the green paradox. *Journal of Environmental Economics and Management*, 68(1):87–110, 2014.
- G. Nemet. Beyond the learning curve: factors influencing cost reductions in photovoltaics. *Energy Policy*, 34(17):3218–3232, 2006.
- W. Nordhaus. *A Question of Balance*. Yale University Press, 2008.
- W. Nordhaus. Estimates of the social cost of carbon: concepts and results from the DICE-2013R model and alternative approaches. *Journal of the Association of Environmental and Resource Economists*, 1(1/2):273–312, 2014a.
- W. Nordhaus. The perils of the learning model for modeling endogenous technological change. *The Energy Journal*, Volume 35(1), 2014b.
- W. Nordhaus. Projections and uncertainties about climate change in an era of minimal climate policies. *American Economic Journal: Economic Policy*, 10(3):333–60, 2018.
- C. Papageorgiou, M. Saam, and P. Schulte. Substitution between clean and dirty energy inputs - a macroeconomic perspective. *Review of Economics and Statistics*, 99(2):281–290, 2017.
- E. Petrakis, E. Rasmusen, and S. Roy. The learning curve in a competitive industry. *The Rand journal of economics*, 28(2):248–268, 1997.
- A. Pfeiffer, R. Millar, C. Hepburn, and E. Beinhocker. The ‘2°C capital stock’ for electricity generation: Committed cumulative carbon emissions from the electricity generation sector and the transition to a green economy. *Applied Energy*, 179:1395–1408, 2016.
- R. S. Pindyck. Irreversibility, uncertainty, and investment. *Journal of Economic Literature*, 29(3):1110–1148, 1991.
- R. S. Pindyck. Pricing capital under mandatory unbundling and facilities sharing. *NBER working paper No. 11225*, 2005.
- R. S. Pindyck. Climate change policy: What do the models tell us? *Journal of Economic Literature*, 51(3):860–72, 2013.
- R. S. Pindyck. The use and misuse of models for climate policy. *Review of Environmental Economics and Policy*, 11(1):100–114, 2017.
- J. Reichenbach and T. Requate. Subsidies for renewable energies in the presence of learning effects and market power. *Resource and Energy Economics*, 34(2):236–254, 2012. ISSN 0928-7655.
- A. Rezai and F. van der Ploeg. Abandoning fossil fuel: How fast and how much. *The Manchester School*, 85:e16–e44, 2017.
- J. Rozenberg, A. Vogt-Schilb, and S. Hallegatte. Instrument choice and stranded assets in the transition to clean capital. *Journal of Environmental Economics and Management*, 2018.
- E. S. Rubin, I. M. Azevedo, P. Jaramillo, and S. Yeh. A review of learning rates for electricity supply technologies. *Energy Policy*, 86:198–218, 2015.
- C. Shearer, N. Ghio, L. Myllyvirta, A. Yu, and T. Nace. Boom and bust 2016: Tracking the global coal plant pipeline. *CoalSwarm, Greenpeace and Sierra Club*, 2016.

- H.-W. Sinn. Public policies against global warming: a supply side approach. *International Tax and Public Finance*, 15 (4):360–394, 2008.
- H.-W. Sinn. Introductory comment; the green paradox: A supply-side view of climate policy. *Review of Environmental Economics and Policy*, 9(2):239–245, 2015.
- S. Smulders, Y. Tsur, and A. Zemel. Announcing climate policy: can a green paradox arise without scarcity? *Journal of Environmental Economics and Management*, 64(3):364–376, 2012.
- A. M. Spence. The learning curve and competition. *The Bell Journal of Economics*, 12(1):49–70, 1981.
- W. Steffen, J. Rockström, K. Richardson, T. M. Lenton, C. Folke, D. Liverman, C. P. Summerhayes, A. D. Barnosky, S. E. Cornell, M. Crucifix, J. F. Donges, I. Fetzer, S. J. Lade, M. Scheffer, R. Winkelmann, and H. J. Schellnhuber. Trajectories of the earth system in the anthropocene. *Proceedings of the Natural Academy of Sciences USA*, August 2018.
- N. Stern. The structure of economic modeling of the potential impacts of climate change: Grafting gross underestimation of risk onto already narrow science models. *Journal of Economic Literature*, 51(3):838–59, 2013.
- N. Stern. Current climate models are grossly misleading. *Nature (Comment)*, 530(7591):407–409, 2016.
- R. Tol. The optimal timing of greenhouse gas emission abatement, the individual rationality and intergenerational equity. In C. Carraro, editor, *International Environmental Agreements on Climate Change*. Kluwer Academic Publishers, 1999.
- G. Unruh. Escaping carbon lock-in. *Energy Policy*, 30:317–325, 2002.
- F. van der Ploeg. Cumulative carbon emissions and the green paradox. *Annual Review of Resource Economics*, 5(1):281–300, 2013.
- F. van der Ploeg and C. Withagen. Growth, renewables, and the optimal carbon tax. *International Economic Review*, 55(1):283–311, 2014.
- A. Vogt-Schilb, G. Meunier, and S. Hallegatte. When starting with the most expensive option makes sense: Optimal timing, cost and sectoral allocation of abatement investment. *Journal of Environmental Economics and Management*, 88:210–233, 2018.
- M. Weitzman. On modeling and interpreting the economics of catastrophic climate change. *The Review of Economics and Statistics*, 91(1):1–19, 2009.
- M. Weitzman. What is the “damage function” for global warming - and what difference might it make? *Climate Change Economics*, 1:57–69, 2010.
- T. P. Wright. Factors affecting the cost of airplanes. *Journal of the Aeronautical Science*, 3(2):122–128, 1936.

# Online Appendix for “To Build or Not to Build? Capital Stocks and Climate Policy” (For Online-Only Publication)

## A Proofs of Theoretical Results: Simplified Model

To start with, we define:

$$P_t := \sum_{s=1}^{\infty} (1-\delta)^{s-1} \Delta_{t,s} (r_{t+s} - \delta - e_{t+s}).$$

This is the net present value of investment in the irreversible asset, infrastructure, relative to the opportunity cost. The following technical lemma is very illuminating:

**Lemma A.1.** *Given the framework above,*

1.  $P_t \leq 0$  for all  $t$ .
2.  $i_t > 0$  only if  $P_t = 0$ .
3.  $i_t > 0$  only if both  $r_t - \delta \leq e_t$  and  $r_{t+1} - \delta \geq e_{t+1}$ .
4.  $i_t > 0$  with  $r_{t+1} - \delta > e_{t+1}$  only if  $i_{t+1} = 0$ .

**Proof of Lemma A.1.** Write  $o_t$  for all other sources of income, net of any other investments (which may also be irreversible). We maximize

$$\sum_{t=1}^{\infty} \beta^t u(c_t)$$

subject to constraints

$$\begin{array}{ll} \mu_t^{bc} & i_t + c_t = r_t k_t + o_t \\ \mu_t^i & i_t \geq 0 \\ \mu_t^k & i_t \geq k_{t+1} - (1-\delta)k_t \end{array}$$

The Lagrangian is:

$$\begin{aligned} \mathcal{L}_t = \sum_{t=0}^{\infty} \beta^t & \left( u(c_t) - \mu_t^{bc}(i_t + c_t) + \mu_t^{bc}(r_t k_t + o_t) + \mu_t^i i_t \right. \\ & \left. + \mu_t^k (i_t - (k_{t+1} - (1-\delta)k_t)) \right) \end{aligned}$$

Leading to FOCs and complementary slack conditions

$$c_t \quad u'(c_t) = \mu_t^{bc} \quad (\text{A.1})$$

$$i_t \quad \mu_t^{bc} = \mu_t^i + \mu_t^k \quad (\text{A.2})$$

$$k_{t+1} \quad \mu_t^k = \beta(\mu_{t+1}^{bc} r_{t+1} + \mu_{t+1}^k (1 - \delta)) \quad (\text{A.3})$$

$$\mu_t^i \geq 0$$

$$\mu_t^i i_t = 0 \quad (\text{A.4})$$

$$\mu_t^k \geq 0 \quad (\text{A.5})$$

$$\mu_t^k (i_t - (k_{t+1} - (1 - \delta)k_t)) = 0$$

Substitute (A.1) and (A.2) into (A.3) and divide by  $\beta u'(c_{t+1})$ :

$$\frac{u'(c_t)}{\beta u'(c_{t+1})} \left(1 - \frac{\mu_t^i}{\mu_t^{bc}}\right) = r_{t+1} + \left(1 - \frac{\mu_{t+1}^i}{\mu_{t+1}^{bc}}\right) (1 - \delta)$$

Write  $e_{t+1} := \frac{u'(c_t)}{\beta u'(c_{t+1})} - 1$  and  $\Delta_{t,s} = \prod_{s'=1}^s \frac{1}{1+e_{t+s'}}$ . Re-arrange so that this will provide a forward-looking formula for  $\frac{\mu_t^i}{\mu_t^{bc}}$ :

$$\begin{aligned} \frac{\mu_t^i}{\mu_t^{bc}} &= \frac{e_{t+1} - (r_{t+1} - \delta)}{e_{t+1} + 1} + \frac{(1 - \delta)}{(e_{t+1} + 1)} \frac{\mu_{t+1}^i}{\mu_{t+1}^{bc}} \\ &= \frac{e_{t+1} - (r_{t+1} - \delta)}{e_{t+1} + 1} + \frac{1 - \delta}{e_{t+1} + 1} \left( \frac{e_{t+2} - (r_{t+2} - \delta)}{e_{t+2} + 1} + \frac{(1 - \delta)}{(e_{t+2} + 1)} \frac{\mu_{t+2}^i}{\mu_{t+2}^{bc}} \right) \\ &= \sum_{s=1}^T (1 - \delta)^{s-1} \Delta_{t,s} (e_{t+s} - r_{t+s} + \delta) + (1 - \delta)^T \Delta_{t,T} \frac{\mu_{t+T}^i}{\mu_{t+T}^{bc}}. \end{aligned} \quad (\text{A.6})$$

Next we will show that the final term in (A.6) tends to zero as  $T \rightarrow \infty$ . Since we assumed that there exist  $\epsilon > 0$  and  $R \gg 0$  with  $-\delta + \epsilon < e_t < R$  for all  $t$ , it follows that  $\frac{1-\delta}{1+e_t} < 1 - \frac{\epsilon}{1+e_t} < 1 - \frac{\epsilon}{R+1}$  for all  $t$  and so that  $(1 - \delta)^T \Delta_{t,T} \rightarrow 0$  as  $T \rightarrow \infty$ . Finally,  $0 \leq \mu_T^i \leq \mu_T^{bc}$  for all  $T$ , by consideration of (A.2) and (A.5). It follows that  $0 \leq \frac{\mu_T^i}{\mu_T^{bc}} \leq 1$ , and hence the final term in (A.6) tends to 0 as  $T \rightarrow \infty$ , and we conclude:

$$\frac{\mu_t^i}{\mu_t^{bc}} = \sum_{s=1}^{\infty} (1 - \delta)^{s-1} \Delta_{t,s} (e_{t+s} - r_{t+s} + \delta) =: -P_t \quad (\text{A.7})$$

with per-period equation: 
$$\Delta_{t,1}^{-1} \frac{\mu_t^i}{\mu_t^{bc}} = (e_{t+1} - r_{t+1} + \delta) + (1 - \delta) \frac{\mu_{t+1}^i}{\mu_{t+1}^{bc}} \quad (\text{A.8})$$

Part 1 of Lemma A.1 follows from (A.7). Next, if  $i_t > 0$ , complementary slackness (A.4) tells us  $\mu_t^i = 0$  and so Part 2 follows from (A.7).

If  $i_t > 0$  then by (A.4)  $\mu_t^i = 0$ , and since  $\mu_{t+1}^i \geq 0$  and  $\mu_{t-1}^i \geq 0$ , (A.8) implies  $r_{t+1} - \delta \geq e_{t+1}$  and  $r_t - \delta \leq e_t$ . In addition, Part 4 follows, in the same way as the previous result: if  $i_t > 0$  with  $r_{t+1} - \delta > e_{t+1}$ , then (A.8) implies  $\mu_{t+1}^i > 0$  and then  $i_{t+1} = 0$  from (A.4).  $\square$

**Proof of Lemma 3.1.** Immediate from Lemma A.1 Part 3.  $\square$

**Proof of Lemma 3.2.** If  $r_{s_1} - \delta < e_{s_1}$  then  $i_{s_1-1} = 0$  (by Lemma A.1 Part 3). However, by assumption,  $i_0 > 0$ . Let  $t_0$  be maximal such that  $t_0 < s_1$  and  $i_{t_0} > 0$ . Now, by Lemma A.1 Part 2,  $P_{t_0} = 0$ . So:

$$\begin{aligned} 0 = P_{t_0} &= \sum_{s=1}^{s_1-t_0} (1-\delta)^{s-1} \Delta_{t_0,s}(r_{t_0+s} - \delta - e_{t_0+s}) + \sum_{s=s_1-t_0+1}^{\infty} (1-\delta)^{s-1} \Delta_{t_0,s}(r_{t_0+s} - \delta - e_{t_0+s}) \\ &= \sum_{s=t_0+1}^{s_1} (1-\delta)^{s-t_0-1} \Delta_{t_0,s-t_0}(r_s - \delta - e_s) + \sum_{s=1}^{\infty} (1-\delta)^{s_1-t_0+s-1} \Delta_{t_0,s_1-t_0+s}(r_{s_1+s} - \delta - e_{s_1+s}) \end{aligned} \quad (\text{A.9})$$

It is easy to show that, for any  $t_1, t_2$ , we have  $\Delta_{0,t_1} \Delta_{t_1,t_2} = \Delta_{0,t_1+t_2}$ . Thus,  $\Delta_{0,t_0} \Delta_{t_0,s_1-t_0+s} = \Delta_{0,s_1+s}$ . It also follows that  $\Delta_{0,s_1} \Delta_{s_1,s} = \Delta_{0,s_1+s}$ , and that  $\Delta_{0,t_0} \Delta_{t_0,s_1-t_0} = \Delta_{0,s_1}$ . Putting these facts together we see that  $\Delta_{t_0,s_1-t_0+s} = \Delta_{t_0,s_1-t_0} \Delta_{s_1,s}$ . So, continuing from (A.9), we see

$$\begin{aligned} P_{t_0} &= \sum_{s=t_0+1}^{s_1} (1-\delta)^{s-t_0-1} \Delta_{t_0,s-t_0}(r_s - \delta - e_s) \\ &\quad + (1-\delta)^{s_1-t_0} \Delta_{t_0,s_1-t_0} \sum_{s=1}^{\infty} (1-\delta)^{s-1} \Delta_{s_1,s}(r_{s_1+s} - \delta - e_{s_1+s}) \\ &= \sum_{s=t_0+1}^{s_1} (1-\delta)^{s-t_0-1} \Delta_{t_0,s-t_0}(r_s - \delta - e_s) + (1-\delta)^{s_1-t_0} \Delta_{t_0,s_1-t_0} P_{s_1}. \end{aligned} \quad (\text{A.10})$$

But  $P_{s_1} \leq 0$  by Lemma A.1 Part 1. And  $i_{t_0} > 0$  so  $r_{t_0} - \delta \leq e_{t_0}$  by Lemma A.1 Part 3. Thus:

$$\sum_{s=t_0+1}^{s_1} (1-\delta)^{s-t_0-1} \Delta_{t_0,s-t_0}(r_s - \delta - e_s) \geq 0.$$

Since  $r_{s_1} - \delta - e_{s_1} < 0$  it follows that there exists  $s \in \{t_0 + 1, \dots, s_1 - 1\}$  such that  $r_s - \delta > e_s$ . Letting  $s_0$  be the minimal such  $s$ , it is clear that this meets our requirements.  $\square$

**Proof of Lemma 3.3.** By exactly the same arguments as those used to prove (A.10), and by  $P_{s_2} \leq 0$ , it follows that

$$\begin{aligned} 0 = P_0 &= \sum_{s=1}^{s_2} (1-\delta)^{s-1} \Delta_{0,s}(r_s - \delta - e_s) + (1-\delta)^{s_2} \Delta_{0,s_2} P_{s_2} \\ &\leq \sum_{s=1}^{s_2} (1-\delta)^{s-1} \Delta_{0,s}(r_s - \delta - e_s) \end{aligned}$$

By splitting the sum into terms with  $s \in \{1, \dots, s_1 - 1\}$  and  $s \in \{s_1, \dots, s_2\}$ , and rearranging, we obtain the expression given.  $\square$

**Proof of Corollary 3.4.** This follows straightforwardly from Lemma 3.2 where we let  $w_2 \rightarrow \infty$ . Let  $t_0$  the minimal time such that  $i_t = 0$  for all  $t \geq t_0$ ; by Lemma 3.2 such a  $t_0$  exists (although it need not be identified with the “ $s_0$ ” found in that Proposition). By definition  $i_{t_0-1} > 0$ . Thus, we



may apply the Proposition using  $t_0 - 1$  as year 0. In particular, then

$$\frac{1}{1-\delta} \sum_{s=t_0}^{s_1} \left( \frac{1-\delta}{1+e} \right)^s ((r_s - \delta) - e) \geq \frac{1}{1-\delta} \sum_{s=s_1}^{\infty} \left( \frac{1-\delta}{1+e} \right)^s (e - (r_s - \delta))$$

Using the assumed bounds for  $r_s - \delta$  in the relevant ranges, it follows that

$$\begin{aligned} \sum_{s=t_0}^{s_1-1} \left( \frac{1-\delta}{1+e} \right)^s d_1 &\geq \sum_{s=s_1}^{\infty} \left( \frac{1-\delta}{1+e} \right)^s d_2 \\ \Rightarrow d_1 \left( 1 - \left( \frac{1-\delta}{1+e} \right)^{s_1-t_0} \right) &\geq d_2 \left( \frac{1-\delta}{1+e} \right)^{s_1-t_0} \\ \Rightarrow \left( \frac{1-\delta}{1+e} \right)^{s_1-t_0} &\leq \frac{d_1}{d_1 + d_2} \\ \Rightarrow s_1 - t_0 &\geq \frac{\log \left( \frac{d_1}{d_1 + d_2} \right)}{\log \left( \frac{1-\delta}{1+e} \right)} = \frac{\log(d_1 + d_2) - \log(d_1)}{\log(1+e) - \log(1+\delta)} \end{aligned}$$

where for clarity we write the numerator and denominator as positive numbers.  $\square$

**Proof of Corollary 3.5.** First, see that without the constraint  $I_t \geq 0$  we have  $\tilde{r}_t - \delta = e_t$  for all  $t$ .

Next, since  $I_0 > 0$  we know  $P_0 = 0$  by Lemma A.1 Part 2. If  $r_t - \delta = e_t = \tilde{r}_t - \delta$  for all  $t$  then  $K_t = \tilde{K}_t$  for all  $t$ , but this is not possible since  $\tilde{I}_{t_1} < 0$  and  $I_{t_1} \geq 0$ . If we assume  $r_t - \delta \geq e_t$  for all  $t$  we must conclude also  $r_t - \delta > e_t$  for some  $t$ , whence  $P_0 > 0$ , which is a contradiction. So there exist some minimal  $s_1$  such that  $r_{s_1} - \delta < e_{s_1}$  and some maximal  $s_2 \in \mathbb{R} \cup \{\infty\}$  such that  $s_2 \geq s_1$  and  $r_t - \delta < e_t$  for  $t \in \{s_1, \dots, s_2\}$ . Applying Lemma 3.2 we conclude that there exists  $s_0 \leq s_1 - 1$  such that  $r_{s_0} - \delta > e_{s_0}$  and such that  $I_t = 0$  for  $t \in \{s_0, \dots, s_2 - 1\}$ . Pick  $s_0$  minimal with these properties.

We show that  $s_0$  is minimal such that  $r_t - \delta \neq e_t$ . First, by definition of  $s_1$ , there is no  $t < s_0$  with  $r_t - \delta < e_t$ . Next, if  $r_t - \delta > e_t$  for  $t < s_0$  then there exists  $t' \in \{t, \dots, s_0 - 1\}$  such that  $I_{t'} > 0$  (for otherwise  $s_0$  is not minimal as defined). But  $P_0 = 0$  and  $P_{t'} = 0$  imply that there must also exist  $t'' \in \{1, \dots, t'\}$  such that  $r_{t''} - \delta < e_{t''}$ , and we already know this is not so.

Since  $r_t - \delta = e_t = \tilde{r}_t - \delta$  for  $t \in \{0, \dots, s_0 - 1\}$ , it follows that  $K_t = \tilde{K}_t$  for  $t \in \{0, \dots, s_0 - 1\}$  and so that  $I_{t-1} = \tilde{I}_{t-1} \geq 0$  for  $t \in \{0, \dots, s_0 - 1\}$ . So we know  $t_1 \geq s_0$ .

Next,  $r_{s_0} - \delta > e_{s_0} = \tilde{r}_{s_0} - \delta$  so  $K_{s_0} < \tilde{K}_{s_0}$ ; but  $K_{s_0-1} = \tilde{K}_{s_0-1}$ , so  $I_{s_0-1} < \tilde{I}_{s_0-1}$ . So set  $t_0 := s_0 - 1$ .

Finally, by definition  $r_{s_1} - \delta < e_{s_1} = \tilde{r}_{s_1} - \delta$ , which implies  $K_{s_1} > \tilde{K}_{s_1}$ . But  $K_{t_0+1} < \tilde{K}_{t_0+1}$  and so, since  $I_t = 0$  for  $t \in \{t_0 + 1, \dots, s_2 - 1\}$  we conclude that  $K_t < \tilde{K}_t$  for  $t \leq \{t_0 + 1, \dots, \min(s_2 - 1, t_1)\}$ . Since  $s_1 \leq s_2 - 1$  and since  $K_{s_1} > \tilde{K}_{s_1}$  we conclude that  $\min(s_2 - 1, t_1) = t_1$ , i.e. that  $K_t < \tilde{K}_t$  for  $t \in \{t_0 + 1, \dots, t_1\}$  as required.  $\square$

**The Social Planner's problem for Section 4.1** The planner maximizes

$$\sum_{t=0}^{\infty} \beta^t L_t u \left( \frac{C_t}{L_t} \right)$$

subject to the constraints:

$$\Lambda_t^s \quad I_t + C_t = f_t(H_t, O_t) \quad (\text{A.11})$$

$$\mu_t^I \quad I_t \geq 0 \quad (\text{A.12})$$

$$\mu_t^H \quad I_t = p_t^H(H_{t+1} - (1 - \delta)H_t) \quad (\text{A.13})$$

$$\mu_t^p \quad p_t^H = G(H_t) \quad (\text{A.14})$$

where  $O_t = L_t o_t$  represents all other factors of production in the economy. In our model the planner treats this as exogenous.

At time  $t$ , the Lagrangian is

$$\begin{aligned} \mathcal{L}_t = \sum_{t=0}^{\infty} \beta^t & \left( L_t u \left( \frac{C_t}{L_t} \right) - \Lambda_t^s (I_t + C_t - f_t(H_t, O_t)) + \mu_t^I I_t \right. \\ & \left. + \mu_t^H (I_t - p_t^H (H_{t+1} - (1 - \delta)H_t)) + \mu_t^p (p_t^H - G(H_t)) \right) \end{aligned}$$

the first order conditions are:

$$\partial C_t : \quad \Lambda_t^s = u' \left( \frac{C_t}{L_t} \right) \quad (\text{A.15})$$

$$\partial H_{t+1} : \quad p_t^H \mu_t^H = \beta \left( \Lambda_{t+1}^s \frac{\partial f_{t+1}}{\partial H_{t+1}} + \mu_{t+1}^H p_{t+1}^H (1 - \delta) \right) - \beta \mu_{t+1}^p G'(H_{t+1}) \quad (\text{A.16})$$

$$\partial I_t : \quad \Lambda_t^s = \mu_t^H + \mu_t^I \quad (\text{A.17})$$

$$\partial p_t^H : \quad \mu_t^p = \mu_t^H (H_{t+1} - (1 - \delta)H_t) \quad (\text{A.18})$$

together with the constraints above and the inequality  $\mu_t^I \geq 0$ , which is complementary slack with (A.12).

**Proof of Proposition 4.1.** Divide (A.16) through by  $p_t^H \beta \Lambda_{t+1}^s$ , substitute in (A.17) and (A.18) and re-arrange to obtain:

$$\begin{aligned} R_{t+1} &= \frac{\mu_t^H - \beta(1 - \delta)\mu_{t+1}^H}{\beta \Lambda_{t+1}^s} \\ &= \frac{1}{p_t^H} \frac{\partial f_{t+1}}{\partial H_{t+1}} + \left( 1 - \frac{\mu_{t+1}^I}{\Lambda_{t+1}^s} \right) \frac{p_{t+1}^H - p_t^H}{p_t^H} (1 - \delta) - \left( 1 - \frac{\mu_{t+1}^I}{\Lambda_{t+1}^s} \right) \frac{H_{t+2} - (1 - \delta)H_{t+1}}{p_t^H} G'(H_{t+1}) \end{aligned}$$

if  $I_{t+1} > 0$  then, by complementary slackness,  $\mu_{t+1}^I = 0$ . Thus, multiplying both sides by  $\frac{p_t^H}{p_{t+1}^H}$ , and substituting in the definition for direct returns we obtain the expression given.

Finally, in the case  $I_{t+1} > 0$ , the defining formula for  $R_{t+1}$  becomes just:

$$R_{t+1} = \frac{\Lambda_t^s - \beta(1 - \delta)\Lambda_{t+1}^s}{\beta \Lambda_{t+1}^s} = \frac{u'(C_t/L_t)}{\beta u'(C_t/L_t)} - 1 + \delta = e_{t+1} + \delta.$$

(where we substitute also from (A.15)), as required.  $\square$

**Corollary A.2.** Assume that  $G'(H) < 0$  and  $H_{t+1} > (1 - \delta)H_t$ . If  $\delta = 1$  then  $\frac{p_t^H}{p_{t+1}^H} R_{t+1} > r_{t+1}^s$ . If  $\delta = 0$ , if  $G$  is convex, and if  $H_{t+2} - H_{t+1} = H_{t+1} - H_t$  then  $\frac{p_t^H}{p_{t+1}^H} R_{t+1} < r_{t+1}^s$ .

Regarding the case in which the price effect dominates the learning effect: the assumption of convexity for the function  $G$  giving the decline in prices, is very natural. The assumption that capital is increasing by a constant *amount*, rather than a constant *factor*, is less so; and increases in  $H_{t+2} - H_{t+1}$  relative to  $H_{t+1} - H_t$  will increase  $\frac{p_t^H}{p_{t+1}^H} R_{t+1}$  relative to  $r_{t+1}^s$ . In general we expect the learning effect to dominate, but it is worth noting that when capital is very persistent, and when there is a considerably delay in realizing the benefits of learning, then the extent to which the learning effect pushes the shadow return above the direct return, is mitigated by the price effect.

**Proof of Proposition 4.2.** Considering first the firm, there is no inter-temporal element to their objective function or constraints and so we can consider their optimization period-by-period; obviously the relevant first-order condition is that

$$\frac{\partial f_t}{\partial H_t} = r_t p_t^H. \quad (\text{A.19})$$

Meanwhile, the household maximizes:

$$\sum_{t=0}^{\infty} \beta^t \frac{L_t}{L_0} u\left(\frac{L_0}{L_t} c_t\right)$$

subject to the constraints:

$$\begin{aligned} \Lambda_t & i_t + c_t = (r_t + \tau_t) p_t^H h_t + o_t \\ \mu_t^i & i_t \geq 0 \\ \mu_t^h & i_t = p_t^H (h_{t+1} - (1 - \delta)h_t) \end{aligned} \quad (\text{A.20})$$

Additionally, the price is constrained by  $p_t^H = G(H_t)$ , but the household does not take this into account. At time  $t$ , the Lagrangian is

$$\begin{aligned} \mathcal{L}_t = & \sum_{t=0}^{\infty} \beta^t \left( \frac{L_t}{L_0} u\left(\frac{L_0}{L_t} c_t\right) - \Lambda_t (i_t + c_t - (r_t + \tau_t) p_t^H h_t - o_t) + \mu_t^i i_t \right. \\ & \left. + \mu_t^h (i_t - p_t^H (h_{t+1} - (1 - \delta)h_t)) \right) \end{aligned}$$

the first order conditions are:

$$\partial c_t : \quad \Lambda_t = u' \left( \frac{L_0}{L_t} c_t \right) = u' \left( \frac{C_t}{L_t} \right) \quad (\text{A.21})$$

$$\partial h_{t+1} : \quad p_t^H \mu_t^h = \beta \Lambda_{t+1} (r_{t+1} + \tau_{t+1}) p_{t+1}^H + \beta \mu_{t+1}^h p_{t+1}^H (1 - \delta) \quad (\text{A.22})$$

$$\partial i_t : \quad \Lambda_t = \mu_t^h + \mu_t^i \quad (\text{A.23})$$

together with the constraints above and the inequality  $\mu_t^i \geq 0$ , which is complementary slack with (A.20).

Substitute (A.19) into (A.22) and rearrange: now this first order condition reads:

$$p_t^H \mu_t^h = \beta \left( \Lambda_{t+1} \frac{\partial f_{t+1}}{\partial H_{t+1}} + \mu_{t+1}^h p_{t+1}^H (1 - \delta) \right) + \beta \Lambda_{t+1} \tau_{t+1} p_{t+1}^H \quad (\text{A.24})$$

We seek the equation for  $\tau_{t+1}$  that will lead to the same solution as in the social planner's problem; as derived above, this is defined by constraints (A.11)–(A.14), first order conditions (A.15)–(A.18) and the inequality  $\mu_t^I \geq 0$ , which is complementary slack with (A.12). Those equations are all counterparts to the equations of this model, with the exception of (A.24): we wish this to imply (A.16). But this will be the case if we set (substituting in also (A.18))

$$\begin{aligned} \Lambda_{t+1} \tau_{t+1} p_{t+1}^H &= -\mu_{t+1}^h (H_{t+2} - (1 - \delta)H_{t+1}) G'(H_{t+1}) \\ \Leftrightarrow \tau_{t+1} &= -\frac{\mu_{t+1}^h}{\Lambda_{t+1}} \frac{H_{t+2} - (1 - \delta)H_{t+1}}{p_{t+1}^H} G'(H_{t+1}) \end{aligned}$$

So if  $i_t > 0$ , which implies  $\mu_{t+1}^h = \Lambda_{t+1}$ , then the two models are defined by the same first-order conditions in variables  $C_t$ ,  $H_t$  and  $I_t$ . In each case  $p_t^H$  is defined by  $H_t$ , so if  $O_t = L_0 o_t$  for all  $t$  then the solutions are equal – that is, this level of subsidy achieves the social optimum (subject to  $O_t$ ).

We have treated  $o_t$  and  $O_t$  as exogenous for both the household and the social planner. More generally, a model will allow optimization in all factors of production and sources of income. However, if all externalities except for the learning-by-doing in  $p_t^H$  have been internalized, then by the Coase Theorem and the First Welfare Theorem, it follows that the optimal  $O_t^*$  for the planner satisfies  $O_t^* = L_0 o_t^*$ , where  $o_t^*$  is optimal for the household, so the solutions to the models coincide. It is straightforward to now write this in terms of  $g_t^H$ .  $\square$

**Proof of Corollary 4.3.** If  $p_t^H = G(H_t) = p_0^H \left( \frac{H_t}{H_0} \right)^{-\lambda}$ , then

$$G'(H_t) = -\lambda \frac{p_0^H}{H_0} \left( \frac{H_t}{H_0} \right)^{-\lambda-1} = -\lambda \frac{p_0^H}{H_t} \left( \frac{H_t}{H_0} \right)^{-\lambda} = -\frac{\lambda p_t^H}{H_t}$$

Hence, in this case,  $\tau_t = \lambda (g_t^H + \delta)$ , as required.  $\square$

## B Calibration

This section describes calibration of the model. We build on the seminal “DICE 2013” climate-economy model of Nordhaus (2014a), which serves as benchmark in the literature and policy applications. Some of the parameter values are drawn from the existing studies, in particular, from Hassler et al. (2012), Papageorgiou et al. (2017) and Rezai and van der Ploeg (2017). All the parameter values are summarized in Table A.1. Details of the calibration are as follows:

### B.1 Production

Labor  $L_0$  is given for 2012 using United Nations data. We assume it continues to evolve as in DICE 2013. We set the value of elasticity of substitution between general output,  $Y^g$ , and electricity,  $E$ , in the final-goods production function,  $\kappa = 0.46$ , following Rezai and van der Ploeg (2017), as a compromise between short-term insubstitutability (Hassler et al. (2012)) and longer-term substitutability. We take the value of  $\theta$  from Papageorgiou et al. (2017) to be 0.003. The technology

weightings  $A_0^g$  and  $A_t^E$  will be set to match other data. Subsequently,  $A_t^g$  evolves as in DICE, and  $A_t^E$  evolves in step with it. We set  $\alpha = 0.4$  as an approximation of the values Papageorgiou et al. (2017) get in their various specifications, but this is also commonly-used value in the literature. We set the depreciation rate of the general capital stock at  $\delta^g = 0.05$  following Rezai and van der Ploeg (2017).

In modeling the electricity sector we follow Papageorgiou et al. (2017): we set the value of  $w$  at 0.32 (across various specifications, they find  $w = 0.19$  to  $0.70$ , with a mean of 0.32). We set the value of the substitution parameter  $\xi = 0.46$ , in line with their estimates. We find the initial generating capital stocks for the dirty and renewable generating capacity from EIA data.<sup>28</sup> We set  $A_0^E$  so that electricity output in the first period matches the EIA data on electricity output in 2012.

In calibrating the prices of fossil and renewable energy capital  $p_t^D, p_t^H$ , we set  $p_t^D$  to be constant and to match the current price of new coal-fired power stations in China, as these may be the marginal new plants in consideration.<sup>29</sup> For  $p_t^H$ , see the section below. Exponential depreciation for fossil fuel and renewable energy capital is calculated so that the net lifetime availability of capital is equal to the general expected lifetime of plants in this sector: 40 and 25 years respectively.

We know the initial value of  $K_t^D$  from EIA data for 2012, and  $D_t$  from European Union data. We assume that initially  $\zeta_t = 1$ .

The function form of fossil fuel extraction cost is taken from Rezai and van der Ploeg (2017), but we calibrate it differently because we are more concerned with the price of coal than oil. So we set  $\gamma_1$  to represent the cost of coal in 2012 (IEA2014 data), which we have converted to give this cost as a price per GtC of CO<sub>2</sub> pollution (so that fuel and pollution will be in a straightforward 1:1 ratio), to give a cost of 0.09 trillion 2010\$ / GtC. We take  $S_0 = 2000$ .<sup>30</sup> Using the IEA estimate of the cost of coal in 2040 along a given trajectory, and the additional fractional fossil stock use that this would represent, the second parameter of the resource cost equation is calculated to be  $\gamma_2 = 1.64$ .

We set the value of  $\phi_2$  in the mitigation expenditure function  $\Psi_t$  from DICE2013.

---

<sup>28</sup>All fossil generating capacity has been included on the “dirty” side. For renewables, we exclude hydropower, because it is a relatively mature source of electricity (costs are not falling very fast) and its use is constrained by physical geography, with a large fraction of suitable sites already in use (so its use cannot expand fast), so this technology does not represent the features of interest in the model. Since extensive hydropower capacity already exists, the inclusion of existing capacity would severely bias the trajectory of the equation relating renewable capital to cost of renewable capital.

<sup>29</sup>Numbers taken from Energy and Environmental Economics, Inc. (2012).

<sup>30</sup>The proven resources of all fossil fuels may be estimated as 1003 GtC using EIA data. However, continued exploration will enlarge these stocks. We use the stock figure of 2000 GtC.

| Parameter           | Value   | Units               | Definition   |
|---------------------|---------|---------------------|--|
| $L_0$               | 7.10    | billion people      | Population   |
| $A_0^g$             | 2.53    |                     | Productivity   |
| $K_0^g$             | 150.00  | trillion 2010\$     | Initial ‘general’ capital stock                          |
| $\theta$            | 0.003   |                     | Energy share parameter, global output                    |
| $\alpha$            | 0.4     |                     | Share of capital, global output                          |
| $\kappa$            | 0.46    |                     | Elas of substitution btw energy and capital/labor        |
| $\xi$               | 0.46    |                     | Elas of subs between clean and dirty electricity capital |
| $\omega$            | 0.32    |                     | Weight on renewable capital in electricity output        |
| $D_0$               | 9.4     | GtC                 | CO <sub>2</sub> Emissions in year 2012                   |
| $D_0^{\text{land}}$ | 0.90    | GtC                 | Land-use CO <sub>2</sub> emissions in year 2012          |
| $D_0^E$             | 3.30    | GtC                 | Electricity CO <sub>2</sub> emissions in year 2012       |
| $D_0^g$             | 5.22    | GtC                 | General economy CO <sub>2</sub> emissions in 2012        |
| $\nu$               | 0.91    | GtC/(tW capacity)   | Fuel use & emissions from dirty electricity production   |
| $S_0$               | 2000    | GtC                 | Existing stock of fossil fuel (as of 2012)               |
| $Y_0$               | 60.11   | trillion 2010\$     | Initial gross world output                               |
| $K_0^D$             | 3.61    | tW                  | Initial capital stock of fossil technology               |
| $H_0$               | 0.46    | tW                  | Initial renewable-energy-knowledge capital stock         |
| $p^D$               | 0.57    | trillion 2010\$/tW  | Price of dirty electricity capital                       |
| $p_0^H$             | 2.11    | trillion 2010\$/tW  | Initial price of clean electricity capital               |
| $\delta^g$          | 0.05    | year <sup>-1</sup>  | Capital stock depreciation rate                          |
| $\delta^D$          | 0.025   | year <sup>-1</sup>  | Fossil energy capital depreciation                       |
| $\delta^H$          | 0.04    | year <sup>-1</sup>  | Renewable energy capital depreciation                    |
| $\gamma_1$          | 0.09    | trillion 2010\$/GtC | Parameter of fuel extraction costs                       |
| $\gamma_2$          | 1.64    |                     | Parameter of fuel extraction costs                       |
| $A_0^E$             | 6.93    |                     | Productivity of energy production                        |
| $\lambda$           | 0.295   |                     | Rate of learning.  |
| $\varsigma_1$       | 0.00267 |                     | Damage function parameter.                               |
| $\varsigma_2$       | 2       |                     | Damage function parameter.                               |
| $\varsigma_3$       | 0.001   |                     | Damage function parameter.                               |
| $\varsigma_4$       | 50      |                     | Damage function parameter.                               |
| $\phi_2$            | 2.8     |                     | Mitigation expenditure parameter.                        |
| $\phi_3$            | 0.01    |                     | Mitigation expenditure parameter.                        |
| $\sigma_0$          | 0.0904  | GtC/trillion 2010\$ | the carbon-equivalent emissions to output ratio.         |
| $\phi_{1,0}$        | 0.041   |                     | Backstop costs.  |

Table A.1: Parameter values

| Variable                          | Definition   |
|-----------------------------------|--|
| $c_t$                             | Per-household consumption  |
| $L_t$                             | Population at period $t$   |
| $K_t^g$                           | Aggregate capital stock in general economy                       |
| $K_t^D$                           | Aggregate dirty capital stock                                    |
| $H_t$                             | Aggregate clean (renewable) capital stock                        |
| $I_t^g$                           | Aggregate investment in general economy                          |
| $I_t^D$                           | Aggregate investment in dirty capital stock                      |
| $I_t^H$                           | Aggregate investment in clean (renewable) capital stock          |
| $\Psi_t$                          | Abatement  |
| $S_t$                             | Fossil fuel stock at period $t$                                  |
| $G^D(S_t)$                        | Fossil fuel extraction costs                                     |
| $r_t^D$                           | Rate of return on fossil (dirty) capital                         |
| $r_t^H$                           | Rate of return on renewable (clean) capital                      |
| $r_t^g$                           | Rate of return on general capital                                |
| $\Pi_t^g$                         | Total profits from sale of the final goods                       |
| $\Pi_t^D$                         | Total profits from sale of the dirty fuel based electricity      |
| $\Pi_t^H$                         | Total profits from sale of the clean electricity                 |
| $\Pi_t^{DE}$                      | Total profits from sale of the fossil fuel                       |
| $\Pi_t^E$                         | Total profits from sale of the aggregate electricity             |
| $\Pi_t$                           | Sum of all profits   |
| $\pi_t$                           | Total profits per-household                                      |
| $p_t^D$                           | Cost of fossil fuel capital                                      |
| $p_t^H$                           | Cost of renewable energy capital                                 |
| $p_t^{EH}$                        | Price of electricity generated by clean power stations           |
| $p_t^{ED}$                        | Price of electricity generated by fossil fuel based power plants |
| $p_t^e$                           | Price of aggregate electricity                                   |
| $p_t^{fuel}$                      | Price of dirty fossil fuel                                       |
| $\Gamma_t^{ED}$                   | Electricity generated by fossil-fuel based power plants          |
| $Y_t = f(Y_t^g, E_t)$             | Total output before damages                                      |
| $Y_t^g$                           | Output of the general economy                                    |
| $E_t = f_t^E(H_t, \Gamma_t^{ED})$ | Aggregate electricity  |
| $\zeta_t$                         | Utilization rate of dirty capital stock                          |
| $\eta_t$                          | Emission control rate in the general sector                      |
| $w_t$                             | Wage   |
| $D_t^E$                           | Fossil fuel (e.g., coal) used in production of electricity       |
| $D_t^g$                           | Fossil fuel used in the general economy                          |

Table A.2: Variables notation and definition

## C The Setup of Social Planner's Problem

We will consider two alternative perspectives for returns on investment, which will be relevant in different contexts. First, as in Section 4.1, we define:

**Definition C.1.** The *shadow returns on investment in the general, dirty and renewable capital*

stocks are defined to be respectively  $R_t^g$ ,  $R_t^D$  and  $R_t^H$  so that:

$$\begin{aligned} R_{t+1}^g &:= \frac{\mu_t^{Kg} - \beta(1 - \delta^g)\mu_{t+1}^{Kg}}{\beta u'(C_{t+1}/L_{t+1})} \\ R_{t+1}^D &:= \frac{\mu_t^{KD} - \beta(1 - \delta^D)\mu_{t+1}^{KD}}{\beta u'(C_{t+1}/L_{t+1})} \\ R_{t+1}^H &:= \frac{\mu_t^{KH} - \beta(1 - \delta^H)\mu_{t+1}^{KH}}{\beta u'(C_{t+1}/L_{t+1})} \end{aligned}$$

where  $\mu_t^{Kg}$ ,  $\mu_t^{KD}$  and  $\mu_t^H$  are the shadow prices on the capital accumulation constraints as below.

On the other hand, one might consider the more immediate definitions for direct economic returns to investment:

**Definition C.2.** The *direct economic returns on investment in the general, dirty and renewable capital stocks* are defined respectively to be  $r_t^g$ ,  $r_t^D$  and  $r_t^H$  so that:

$$\begin{aligned} r_{t+1}^g &:= \frac{\partial}{\partial K_{t+1}^g} (Y_{t+1} - \Psi_{t+1}) \\ r_{t+1}^D &:= \frac{1}{p_{t+1}^D} \frac{\partial}{\partial K_{t+1}^D} (Y_{t+1} - \Psi_{t+1}) \\ r_{t+1}^H &:= \frac{1}{p_{t+1}^H} \frac{\partial}{\partial H_{t+1}} (Y_{t+1} - \Psi_{t+1}) \end{aligned}$$

Here we measure the direct effects of investment on output net of mitigation costs, and the output is

$$Y_t = \Omega(T_t)f(Y_t^g, E_t) \quad (\text{A.25})$$

with  $Y_t^g = f_t^g(K_t^g, L_t)$ .

The social planner's problem is outlined below. Specifically, the social planner maximizes the social welfare function:

$$\sum_{t=0}^{\infty} \beta^t L_t u\left(\frac{C_t}{L_t}\right) \quad (\text{A.26})$$



subject to constraints:

$$Y_t = I_t^g + I_t^D + I_t^H + C_t + G^D(S_t)(D_t^E + D_t^g) + \frac{\phi_{1,t}\eta_t^{\phi_2}Y_t^g}{(1-\eta_t)^{\phi_3}} \quad \mu_t^{BC} \quad (\text{A.27})$$

$$S_{t+1} = S_t - D_t^E - D_t^g \quad \mu_t^S \quad (\text{A.28})$$

$$D_t = D_t^E + D_t^{\text{land}} + D_t^g \quad \mu_t^D \quad (\text{A.29})$$

$$T_t = \mathcal{W}_t(D_0, \dots, D_{t-1}) \quad \mu_t^W \quad (\text{A.30})$$

$$E_t = f_t^E(H_t, \zeta_t K_t^D) = A_t^E \left( \omega(H_t)^\xi + (1-\omega)(\zeta_t K_t^D)^\xi \right)^{1/\xi} \quad \mu_t^E \quad (\text{A.31})$$

$$D_t^E = \nu \zeta_t K_t^D \quad \mu_t^{DE} \quad (\text{A.32})$$

$$D_t^g = \sigma_t(1-\eta_t)Y_t^g \quad \mu_t^{Dg} \quad (\text{A.33})$$

$$\zeta_t \leq 1 \quad \mu_t^\zeta \quad (\text{A.34})$$

$$p_t^H = G(H_t) \quad \mu_t^{pH} \quad (\text{A.35})$$

$$I_t^g = K_{t+1}^g - (1-\delta^g)K_t^g \quad \mu_t^{Kg} \quad (\text{A.36})$$

$$I_t^D = p^D(K_{t+1}^D - (1-\delta^D)K_t^D) \quad \mu_t^{KD} \quad (\text{A.37})$$

$$I_t^H = p_t^H(H_{t+1} - (1-\delta^H)H_t) \quad \mu_t^{KH} \quad (\text{A.38})$$

$$I_t^D \geq 0 \quad \mu_t^{ID} \quad (\text{A.39})$$

$$I_t^H \geq 0 \quad \mu_t^{IH} \quad (\text{A.40})$$

(We do not need to specify  $\zeta_t \geq 0$  as this will never be violated in the optimum.) So we calculate the Lagrangian  $\mathcal{L}$  as

$$\begin{aligned} \mathcal{L} = & \sum_{t=0}^{\infty} \beta^t \left[ L_t u \left( \frac{C_t}{L_t} \right) - \mu_t^S (S_{t+1} - S_t + D_t^E + D_t^g) + \mu_t^D (D_t - D_t^E - D_t^{\text{land}} - D_t^g) \right] \\ & + \sum_{t=0}^{\infty} \beta^t \mu_t^W (T_t - \mathcal{W}_t(D_0, \dots, D_{t-1})) \\ & + \sum_{t=0}^{\infty} \beta^t \mu_t^{BC} \left[ \Omega(T_t) f(Y_t^g, E_t) - I_t^g - I_t^D - I_t^H - C_t - G^D(S_t)(D_t^E + D_t^g) - \frac{\phi_{1,t}\eta_t^{\phi_2}Y_t^g}{(1-\eta_t)^{\phi_3}} \right] \\ & - \sum_{t=0}^{\infty} \beta^t [\mu_t^E (E_t - f_t^E(H_t, \zeta_t K_t^D))] \\ & + \sum_{t=0}^{\infty} \beta^t \left[ \mu_t^{DE} (D_t^E - \nu \zeta_t K_t^D) + \mu_t^{Dg} (D_t^g - \sigma_t(1-\eta_t)Y_t^g) + \mu_t^{pH} (p_t^H - G(H_t)) + \mu_t^\zeta (1 - \zeta_t) \right] \\ & + \sum_{t=0}^{\infty} \beta^t [\mu_t^{Kg} (I_t^g - K_{t+1}^g + (1-\delta^g)K_t^g)] \\ & + \sum_{t=0}^{\infty} \beta^t [\mu_t^{KD} (I_t^D - p^D K_{t+1}^D + p^D(1-\delta^D)K_t^D) + \mu_t^{ID} I_t^D] \\ & + \sum_{t=0}^{\infty} \beta^t [\mu_t^{KH} (I_t^H - p_t^H H_{t+1} + p_t^H(1-\delta^H)H_t) + \mu_t^{IH} I_t^H] \end{aligned} \quad (\text{A.41})$$

We obtain the following first order conditions (using shorthand  $f_t$  for  $f(Y_t^g, E_t)$ ,  $f_t^g$  for  $f_t^g(K_t^g, L_t)$ , etc.)

$$\partial C_t : \quad u' \left( \frac{C_t}{L_t} \right) = \mu_t^{BC} \quad (\text{A.42})$$

$$\partial S_{t+1} : \quad \beta \mu_{t+1}^S = \mu_t^S + \beta \mu_{t+1}^{BC} \frac{dG^D}{dS}(S_{t+1})(D_{t+1}^E + D_{t+1}^g) \quad (\text{A.43})$$

$$\partial D_t^E : \quad \mu_t^{DE} = \mu_t^S + \mu_t^D + \mu_t^{BC} G^D(S_t) \quad (\text{A.44})$$

$$\partial D_t^g : \quad \mu_t^{Dg} = \mu_t^S + \mu_t^D + \mu_t^{BC} G^D(S_t) \quad (\text{A.45})$$

$$\partial D_t : \quad \mu_t^D = \sum_{m=0}^{\infty} \beta^m \mu_{t+m}^W \frac{\partial \mathcal{W}_{t+m}}{\partial D_t} \quad (\text{A.46})$$

$$\partial T_t : \quad \mu_t^W = -\mu_t^{BC} \Omega'(T_t) f_t \quad (\text{A.47})$$

$$\partial E_t : \quad \mu_t^E = \mu_t^{BC} \Omega(T_t) \frac{\partial f_t}{\partial E_t} \quad (\text{A.48})$$

$$\partial K_{t+1}^g : \quad \mu_t^{Kg} = \beta \mu_{t+1}^{Kg} (1 - \delta^g) + \beta \mu_{t+1}^{BC} \left( \Omega(T_{t+1}) \frac{\partial f_{t+1}}{\partial Y_{t+1}^g} - \frac{\phi_{1,t+1} \eta_{t+1}^{\phi_2}}{(1 - \eta_{t+1})^{\phi_3}} \right) \frac{\partial f_{t+1}^g}{\partial K_{t+1}^g} \quad (\text{A.49})$$

$$\partial I_t^g : \quad \mu_t^{Kg} = \mu_t^{BC} \quad (\text{A.50})$$

$$\partial K_{t+1}^D : \quad p^D \mu_t^{KD} = \beta p^D \mu_{t+1}^{KD} (1 - \delta^D) + \beta \zeta_{t+1} \left( \mu_{t+1}^E \frac{\partial f_{t+1}^E}{\partial (\zeta_{t+1} K_{t+1}^D)} - \mu_{t+1}^{DE} \nu \right) \quad (\text{A.51})$$

$$\partial I_t^D : \quad \mu_t^{KD} = \mu_t^{BC} - \mu_t^{ID} \quad (\text{A.52})$$

$$\partial H_{t+1} : \quad p_t^H \mu_t^{KH} = \beta p_{t+1}^H \mu_{t+1}^{KH} (1 - \delta^H) + \beta \mu_{t+1}^E \frac{\partial f_{t+1}^E}{\partial H_{t+1}} - \beta \mu_{t+1}^{pH} G'(H_{t+1}) \quad (\text{A.53})$$

$$\partial I_t^H : \quad \mu_t^{KH} = \mu_t^{BC} - \mu_t^{IH} \quad (\text{A.54})$$

$$\partial p_t^H : \quad \mu_t^{pH} = \mu_t^{KH} (H_{t+1} - (1 - \delta^H) H_t) \quad (\text{A.55})$$

$$\partial \zeta : \quad \mu_t^\zeta = K_t^D \left( \mu_t^E \frac{\partial f_t^E}{\partial (\zeta_t K_t^D)} - \mu_t^{DE} \nu \right) \quad (\text{A.56})$$

$$\partial \eta_t : \quad \sigma_t \mu_t^{Dg} = \mu_t^{BC} \frac{\phi_{1,t} \eta_t^{\phi_2 - 1}}{(1 - \eta_t)^{1 + \phi_3}} [\phi_2 (1 - \eta_t) + \eta_t \phi_3] \quad (\text{A.57})$$

together with constraints (A.27)–(A.40) and inequalities  $\mu_t^\zeta \geq 0$ ,  $\mu_t^{ID} \geq 0$ ,  $\mu_t^{IH} \geq 0$  which are complementary slack with corresponding equations (A.34) and (A.39)–(A.40).

It is useful to prove the following proposition, which gives an expression for the rates of return plotted in Figure 2.

**Proposition C.3.** *In the optimal social planner's solution,*

$$\begin{aligned} R_{t+1}^g &= r_{t+1}^g \\ R_{t+1}^D &= \frac{\zeta_t}{p^D} \left( \frac{\partial Y_{t+1}}{\partial (\zeta_t K_{t+1}^D)} - \nu \frac{\mu_{t+1}^{DE}}{u'(C_{t+1}/L_{t+1})} \right) = r_{t+1}^D - \frac{\nu \zeta_t}{p^D} \frac{\mu_{t+1}^{DE}}{u'(C_{t+1}/L_{t+1})} \end{aligned}$$

where  $\mu_{t+1}^{DE}$  is the shadow price on constraint (A.32), which determines emissions  $D_{t+1}^E$  from the use of dirty energy capital.

Returns on general capital take into account the mitigation expense involved in the use of general capital. Returns on dirty energy capital are given by the marginal productivity of fossil fuel capital infrastructure in output, net of the shadow price of fuel used alongside it, and scaled by the price of this capital stock. Considering what Definition C.2 and Proposition C.3 tell us about general capital, throughout the rest of the paper, we refer to  $r_t^g$  as *the* return on investment in the general capital stock.

**Proof of Proposition C.3.** Substitute (A.42) and (A.50) into (A.49) and use Definition C.1 to prove that  $R_{t+1}^g = r_{t+1}^g$ . It has been presented in a more compact form from the observations that  $Y_{t+1} = Z(T_{t+1})f_{t+1}$  and  $\Psi_{t+1} = \frac{\phi_{1,t+1}\eta_{t+1}^{\phi_2}}{(1-\eta_{t+1})^{\phi_3}}Y_{t+1}^g$ . The form that is most useful for further derivations is (from (A.50)):

$$R_{t+1}^g - \delta^g = \frac{\mu_t^{BC}}{\beta\mu_{t+1}^{BC}} - 1. \quad (\text{A.58})$$

For  $R_{t+1}^D$ , divide (A.51) by  $\beta p^D \mu_{t+1}^{BC}$  and substitute (A.48), and then (A.52) and (A.42) to obtain

$$\begin{aligned} \frac{\mu_t^{KD}}{\beta\mu_{t+1}^{BC}} &= \frac{\mu_{t+1}^{KD}}{\mu_{t+1}^{BC}}(1 - \delta^D) + \frac{\zeta_{t+1}}{p^D} \left( Z(T_{t+1}) \frac{\partial f_{t+1}}{\partial E_{t+1}} \frac{\partial f_{t+1}^E}{\partial(\zeta_{t+1}K_{t+1}^D)} - \frac{\mu_{t+1}^{DE}}{\mu_{t+1}^{BC}} \nu \right) \\ \Rightarrow \frac{\zeta_{t+1}}{p^D} \left( \frac{\partial Y_{t+1}}{\partial(\zeta_{t+1}K_{t+1}^D)} - \frac{\mu_{t+1}^{DE}}{\mu_{t+1}^{BC}} \nu \right) &= \frac{\mu_t^{KD}}{\beta\mu_{t+1}^{BC}} - \frac{\mu_{t+1}^{KD}}{\mu_{t+1}^{BC}}(1 - \delta^D) = R_{t+1}^D \end{aligned}$$

as required.  $\square$

To prove Proposition 5.2 of the main text, we will use the following results.

**Proposition C.4.** [*The social cost of carbon*] In an optimal solution:

$$\chi_t = -u' \left( \frac{C_t}{L_t} \right)^{-1} \sum_{m=1}^{\infty} \beta^m u' \left( \frac{C_{t+m}}{L_{t+m}} \right) \frac{\partial Y_{t+m}}{\partial D_t}. \quad (\text{A.59})$$

**Proof of Proposition C.4 (Social Cost of Carbon).** Substitute (A.47) into (A.46), and divide through by  $\mu_t^{BC}$ , to obtain:

$$\begin{aligned} \frac{\mu_t^D}{\mu_t^{BC}} &= - \sum_{m=0}^{\infty} \beta^m \left( \frac{\mu_{t+m}^{BC}}{\mu_t^{BC}} \Omega'(T_{t+m}) f_{t+m} \right) \frac{\partial \mathcal{W}_{t+m}}{\partial D_t} \\ &= - \sum_{m=1}^{\infty} \beta^m \left( \frac{\mu_{t+m}^{BC}}{\mu_t^{BC}} \Omega'(T_{t+m}) f_{t+m} \right) \frac{\partial \mathcal{W}_{t+m}}{\partial D_t} \end{aligned} \quad (\text{A.60})$$

where the sum is from  $m = 1$  because  $\frac{\partial \mathcal{W}_t}{\partial D_t} = 0$ . Next, note that

$$\frac{\partial Y_{t+m}}{\partial D_t} = \Omega'(T_{t+m}) \frac{\partial \mathcal{W}_{t+m}}{\partial D_t} f_{t+m} \quad (\text{A.61})$$

Substituting (A.61), as well as (A.42), into (A.60), we obtain and write this as

$$\chi_t := \frac{\mu_t^D}{u'(C_t/L_t)} = \frac{\mu_t^D}{\mu_t^{BC}} = -u' \left( \frac{C_t}{L_t} \right)^{-1} \sum_{m=1}^{\infty} \beta^m u' \left( \frac{C_{t+m}}{L_{t+m}} \right) \frac{\partial Y_{t+m}}{\partial D_t}. \quad (\text{A.62})$$

Since  $\Omega'(T_{t+m}) < 0$ , we have  $\partial Y_{t+m}/\partial D_t < 0$ , then  $\chi_t > 0$ . We call this term the social cost of carbon (SCC). It represents the marginal future welfare effect of emissions in terms of current welfare.  $\square$

**Proposition C.5. [Hotelling with fossil stocks]** Write  $\mu_t^S$  for the shadow price on Equation (A.28) constraining the stock of fossil fuel. Then:

$$\frac{\mu_{t+1}^S}{u'(C_{t+1}/L_{t+1})} = \frac{\mu_t^S}{u'(C_t/L_t)} (1 - \delta^g + r_{t+1}^g) + \frac{dG^D}{dS}(S_{t+1})(D_{t+1}^E + D_{t+1}^g) \quad (\text{A.63})$$

$$\text{and so } \frac{\mu_t^S}{u'(C_t/L_t)} = - \sum_{s=1}^{\infty} \Delta_{t,s} (G^D)'(S_{t+s})(D_{t+s}^E + D_{t+s}^g) \quad (\text{A.64})$$

where  $\Delta_{t,s} = \prod_{s'=1}^s \frac{1}{1 - \delta^g + r_{t+s'}^g}$  is the compound discount factor.

That is, the return on extracting a unit of fossil fuels tomorrow should be equal to the return on extracting an extra unit today, selling it and getting a return on it at the rate of interest, less the increase in future extraction cost.

**Proof of Proposition C.5 (Hotelling with fossil stocks).** Divide (A.43) through by  $\mu_{t+1}^{BC}$ :

$$\beta \frac{\mu_{t+1}^S}{\mu_{t+1}^{BC}} = \frac{\mu_t^S}{\mu_t^{BC}} \frac{\mu_t^{BC}}{\mu_{t+1}^{BC}} + \beta \frac{dG^D}{dS}(S_{t+1})(D_{t+1}^E + D_{t+1}^g)$$

Substitute in (A.58) and divide by  $\beta$ , to obtain the Hotelling rule:

$$\frac{\mu_{t+1}^S}{\mu_{t+1}^{BC}} = \frac{\mu_t^S}{\mu_t^{BC}} (1 - \delta^g + r_{t+1}^g) + \frac{dG^D}{dS}(S_{t+1})(D_{t+1}^E + D_{t+1}^g)$$

That is, we proved Equation (A.63) as  $\mu_t^{BC} = u'(C_t/L_t)$  from (A.42). To get the infinite sum, repeatedly substitute:

$$\begin{aligned} \frac{\mu_t^S}{\mu_t^{BC}} &= \frac{1}{1 - \delta^g + r_{t+1}^g} \left( \frac{\mu_{t+1}^S}{\mu_{t+1}^{BC}} - (G^D)'(S_{t+1})(D_{t+1}^E + D_{t+1}^g) \right) \\ &= \frac{1}{1 - \delta^g + r_{t+1}^g} \left( \frac{1}{1 - \delta^g + r_{t+2}^g} \left( \frac{\mu_{t+2}^S}{\mu_{t+2}^{BC}} - (G^D)'(S_{t+2})(D_{t+2}^E + D_{t+2}^g) \right) \right. \\ &\quad \left. - (G^D)'(S_{t+1})(D_{t+1}^E + D_{t+1}^g) \right) \\ &= - \sum_{s=1}^{\infty} \Delta_{t,s} (G^D)'(S_{t+s})(D_{t+s}^E + D_{t+s}^g) \end{aligned}$$

where  $\Delta_{t,s} = \prod_{s'=1}^s \frac{1}{1 - \delta^g + r_{t+s'}^g}$ . That is, we have proved Equation (A.64).  $\square$

**Proposition C.6. [Returns on dirty fuel]**

$$\frac{\partial Y_t}{\partial D_t^E} = \frac{\mu_t^S}{u'(C_t/L_t)} + \chi_t + G^D(S_t) + \frac{p^D R_t^D}{\zeta_t \nu}.$$

That is, in an optimal solution, the marginal productivity of fossil fuels in the final output is equal to the shadow value of fossil fuel stocks plus the SCC, the extraction cost, and the fraction of the rate of return on investment in  $K^D$  (gross of depreciation) which represents fuel use.

**Proof of Proposition C.6 (Returns on dirty fuel).** Now take (A.44), divide by  $\mu_t^{BC}$  and substitute in (A.59):

$$\frac{\mu_t^{DE}}{\mu_t^{BC}} = \frac{\mu_t^S}{\mu_t^{BC}} + \chi_t + G^D(S_t)$$

For  $R_{t+1}^D$ , divide (A.51) by  $\beta p^D \mu_{t+1}^{BC}$  and substitute (A.48), and then (A.52) and (A.42) to obtain

$$\begin{aligned} \frac{\mu_t^{KD}}{\beta \mu_{t+1}^{BC}} &= \frac{\mu_{t+1}^{KD}}{\mu_{t+1}^{BC}} (1 - \delta^D) + \frac{\zeta_{t+1}}{p^D} \left( \Omega(T_{t+1}) \frac{\partial f_{t+1}}{\partial E_{t+1}} \frac{\partial f_{t+1}^E}{\partial (\zeta_{t+1} K_{t+1}^D)} - \frac{\mu_{t+1}^{DE}}{\mu_{t+1}^{BC}} \nu \right) \\ \Rightarrow \frac{\zeta_{t+1}}{p^D} \left( \frac{\partial Y_{t+1}}{\partial (\zeta_{t+1} K_{t+1}^D)} - \frac{\mu_{t+1}^{DE}}{\mu_{t+1}^{BC}} \nu \right) &= \frac{\mu_t^{KD}}{\beta \mu_{t+1}^{BC}} - \frac{\mu_{t+1}^{KD}}{\mu_{t+1}^{BC}} (1 - \delta^D) = R_{t+1}^D \end{aligned}$$

which could be written as:

$$R_{t+1}^D = \frac{\zeta_{t+1}}{p^D} \left( \Omega(T_{t+1}) \frac{\partial f_{t+1}}{\partial E_{t+1}} \frac{\partial f_{t+1}^E}{\partial (\zeta_{t+1} K_{t+1}^D)} - \nu \left( \frac{\mu_{t+1}^S}{\mu_{t+1}^{BC}} + \chi_{t+1} + G^D(S_{t+1}) \right) \right).$$

Now, differentiating (A.25) by  $D_t^E$  and multiplying by  $\nu$ :

$$\nu \frac{\partial Y_{t+1}}{\partial D_{t+1}^E} = \Omega(T_{t+1}) \frac{\partial f_{t+1}}{\partial E_{t+1}} \frac{\partial f_{t+1}^E}{\partial (\zeta_{t+1} K_{t+1}^D)}$$

So:

$$\begin{aligned} R_{t+1}^D &= \frac{\nu \zeta_{t+1}}{p^D} \left( \frac{\partial Y_{t+1}}{\partial D_{t+1}^E} - \frac{\mu_{t+1}^S}{\mu_{t+1}^{BC}} - \chi_{t+1} - G^D(S_{t+1}) \right) \\ \Rightarrow \frac{\partial Y_t}{\partial D_t^E} &= \frac{\mu_t^S}{\mu_t^{BC}} + \chi_t + G^D(S_t) + \frac{p^D R_t^D}{\zeta_t \nu} \end{aligned}$$

□

**Lemma C.7.** In the optimal social planner's solution, If  $I_t^H > 0$  and  $I_{t+1}^H > 0$  then:

$$\frac{p_{t+1}^H}{p_t^H} r_{t+1}^H = 1 + r_{t+1}^g - \delta^g - \frac{p_{t+1}^H}{p_t^H} (1 - \delta^H) + \frac{(H_{t+2} - (1 - \delta^H) H_{t+1})}{p_t^H} G'(H_{t+1})$$

**Proof of Lemma C.7.** Consider the equation for renewable capital (A.53). Dividing by  $\beta p_t^H \mu_{t+1}^{BC}$ , and substituting in equations (A.48) and (A.55) as well as (A.54), we see

$$\begin{aligned} \frac{\mu_t^{KH}}{\beta \mu_{t+1}^{BC}} &= \frac{p_{t+1}^H}{p_t^H} \frac{(\mu_{t+1}^{BC} - \mu_{t+1}^{IH})}{\mu_{t+1}^{BC}} (1 - \delta^H) + \frac{\Omega(T_{t+1})}{p_t^H} \frac{\partial f_{t+1}}{\partial E_{t+1}} \frac{\partial f_{t+1}^E}{\partial H_{t+1}} \\ &\quad - \frac{(\mu_{t+1}^{BC} - \mu_{t+1}^{IH})}{\mu_{t+1}^{BC}} \frac{(H_{t+2} - (1 - \delta^H)H_{t+1})}{p_t^H} G'(H_{t+1}). \\ &= \left(1 + \frac{p_{t+1}^H - p_t^H}{p_t^H}\right) \left(1 - \frac{\mu_{t+1}^{IH}}{\mu_{t+1}^{BC}}\right) (1 - \delta^H) + \frac{1}{p_t^H} \frac{\partial Y_{t+1}}{\partial H_{t+1}} \\ &\quad - \left(1 - \frac{\mu_{t+1}^{IH}}{\mu_{t+1}^{BC}}\right) \frac{(H_{t+2} - (1 - \delta^H)H_{t+1})}{p_t^H} G'(H_{t+1}). \end{aligned}$$

From (A.54) and (A.58), we have

$$\frac{\mu_t^{KH}}{\beta \mu_{t+1}^{BC}} = \frac{\mu_t^{BC} - \mu_t^{IH}}{\beta \mu_{t+1}^{BC}} = (1 + r_{t+1}^g - \delta^g) \left(1 - \frac{\mu_t^{IH}}{\mu_t^{BC}}\right)$$

Combining the above two equations and Definition C.2, we have

$$\begin{aligned} (1 + r_{t+1}^g - \delta^g) \left(1 - \frac{\mu_t^{IH}}{\mu_t^{BC}}\right) &= \left(1 + \frac{p_{t+1}^H - p_t^H}{p_t^H}\right) \left(1 - \frac{\mu_{t+1}^{IH}}{\mu_{t+1}^{BC}}\right) (1 - \delta^H) + r_{t+1}^H \frac{p_{t+1}^H}{p_t^H} \\ &\quad - \left(1 - \frac{\mu_{t+1}^{IH}}{\mu_{t+1}^{BC}}\right) \frac{(H_{t+2} - (1 - \delta^H)H_{t+1})}{p_t^H} G'(H_{t+1}). \end{aligned}$$

This gives the more general form; when  $I_t^H > 0$  and  $I_{t+1}^H > 0$ , implying  $\mu_t^{IH} = \mu_{t+1}^{IH} = 0$ , then the version given in the lemma follows.  $\square$

## D Decentralized Equilibrium

A representative household maximizes:

$$\sum_{t=0}^{\infty} \beta^t \frac{L_t}{L_0} u\left(\frac{L_0}{L_t} c_t\right)$$

subject to the constraints:

$$\Lambda_t \quad i_t^g + i_t^D + i_t^H + c_t = \frac{L_t}{L_0} w_t + \pi_t + r_t^g k_t^g + r_t^D p_t^D k_t^D + r_t^H p_t^H h_t + \frac{1}{L_0} (\tau_t^D (D_t^E + D_t^g) - \tau_t^H p_t^H H_t) \quad (\text{A.65})$$

$$\mu_t^{iD} \quad i_t^D \geq 0 \quad (\text{A.66})$$

$$\mu_t^{iH} \quad i_t^H \geq 0$$

$$\mu_t^{kg} \quad i_t^g = k_{t+1}^g - (1 - \delta^g) k_t^g$$

$$\mu_t^{kD} \quad i_t^D = p_t^D (k_{t+1}^D - (1 - \delta^D) k_t^D)$$

$$\mu_t^{kH} \quad i_t^H = p_t^H (k_{t+1}^H - (1 - \delta^H) k_t^H)$$

At time  $t$ , the Lagrangian is

$$\begin{aligned} \mathcal{L}_t = & \sum_{t=0}^{\infty} \beta^t \left( \frac{L_t}{L_0} u \left( \frac{L_0}{L_t} c_t \right) - \Lambda_t (i_t^g + i_t^D + i_t^H + c_t) + \Lambda_t \left( \frac{L_t}{L_0} w_t + \pi_t + r_t^g k_t^g + r_t^D p_t^D k_t^D + r_t^H p_t^H h_t \right) \right. \\ & + \frac{\Lambda_t}{L_0} (\tau_t^D (D_t^E + D_t^g) - \tau_t^H p_t^H H_t) + \mu_t^{iD} i_t^D + \mu_t^{iH} i_t^H + \mu_t^{kg} (i_t^g - (k_{t+1}^g - (1 - \delta^g) k_t^g)) \\ & \left. + \mu_t^{kD} (i_t^D - p_t^D (k_{t+1}^D - (1 - \delta^D) k_t^D)) + \mu_t^{kH} (i_t^H - p_t^H (k_{t+1}^H - (1 - \delta^H) k_t^H)) \right) \end{aligned}$$

the first order conditions are:

$$\partial c_t : \quad \Lambda_t = u' \left( \frac{L_0}{L_t} c_t \right) = u' \left( \frac{C_t}{L_t} \right) \quad (\text{A.67})$$

$$\partial k_{t+1}^g : \quad \mu_t^{kg} = \beta (\Lambda_{t+1} r_{t+1}^g + \mu_{t+1}^{kg} (1 - \delta^g)) \quad (\text{A.68})$$

$$\partial k_{t+1}^D : \quad p_t^D \mu_t^{kD} = \beta (\Lambda_{t+1} p_{t+1}^D r_{t+1}^D + \mu_{t+1}^{kD} p_{t+1}^D (1 - \delta^D)) \quad (\text{A.69})$$

$$\partial h_{t+1} : \quad p_t^H \mu_t^{kH} = \beta (\Lambda_{t+1} p_{t+1}^H r_{t+1}^H + \mu_{t+1}^{kH} p_{t+1}^H (1 - \delta^H)) \quad (\text{A.70})$$

$$\partial i_t^g : \quad \Lambda_t = \mu_t^{kg} \quad (\text{A.71})$$

$$\partial i_t^D : \quad \Lambda_t = \mu_t^{kD} + \mu_t^{iD} \quad (\text{A.72})$$

$$\partial i_t^H : \quad \Lambda_t = \mu_t^{kH} + \mu_t^{iH} \quad (\text{A.73})$$

together with the constraints above and the inequalities  $\mu_t^{iD} \geq 0$ ,  $\mu_t^{iH} \geq 0$ , which are complementary slack with (A.65) and (A.66).

As usual we combine (A.68) with (A.71) to write:

$$\frac{\Lambda_t}{\beta \Lambda_{t+1}} = 1 - \delta^g + r_{t+1}^g \quad (\text{A.74})$$

Substitute (A.73) into (A.70), divide by  $\Lambda_{t+1}$ , and then substitute in (A.74) and divide by  $\beta$ :

$$\begin{aligned}
p_t^H \left( \frac{\Lambda_t}{\Lambda_{t+1}} - \frac{\mu_t^{iH}}{\Lambda_{t+1}} \right) &= \beta \left( p_{t+1}^H r_{t+1}^H + \left( 1 - \frac{\mu_{t+1}^{iH}}{\Lambda_{t+1}} \right) p_{t+1}^H (1 - \delta^H) \right) \\
\Leftrightarrow p_t^H (1 - \delta^g + r_{t+1}^g) \left( 1 - \frac{\mu_t^{iH}}{\Lambda_t} \right) &= p_{t+1}^H r_{t+1}^H + \left( 1 - \frac{\mu_{t+1}^{iH}}{\Lambda_{t+1}} \right) p_{t+1}^H (1 - \delta^H) \\
\Leftrightarrow p_{t+1}^H r_{t+1}^H &= p_t^H (1 - \delta^g + r_{t+1}^g) \left( 1 - \frac{\mu_t^{iH}}{\Lambda_t} \right) - p_{t+1}^H (1 - \delta^H) \left( 1 - \frac{\mu_{t+1}^{iH}}{\Lambda_{t+1}} \right)
\end{aligned}$$

Recall that also  $\mu_t^{iH} i_t^H = 0$ . So we will be able to combine this result with others below to obtain equations determining  $i_t^H$  and, thus, we will be able to scale up the household's problem.

Similarly, considering dirty capital, we can substitute (A.72) into (A.69), then substitute in (A.74) to obtain:

$$p_{t+1}^D r_{t+1}^D = p_t^D (1 - \delta^g + r_{t+1}^g) \left( 1 - \frac{\mu_t^{iD}}{\Lambda_t} \right) - p_{t+1}^D (1 - \delta^D) \left( 1 - \frac{\mu_{t+1}^{iD}}{\Lambda_{t+1}} \right)$$

And, again,  $\mu_t^{iD} i_t^D = 0$ .

Of course, if investment is ongoing ( $\mu_t^{iH} = \mu_{t+1}^{iH} = \mu_t^{iD} = \mu_{t+1}^{iD} = 0$ ) then these two equations are identities between variables we are claiming are “exogenous”. In that case, these provide necessary conditions on investment being non-zero (and non-infinite).

Moreover, because the economy is made up of identical agents behaving in this same way, we may sum complementary slack equations over all these agents to obtain

$$\begin{aligned}
\mu_t^{iH} I_t^H &= 0 \\
\mu_t^{iD} I_t^D &= 0
\end{aligned}$$

Moreover, now we have equations for the solution to the maximization problem, we can scale up from the household level. We have determined that, given prices and rates of return (equations for which follow) aggregate consumption  $C_t$  and investments  $I_t^g$ ,  $I_t^D$ ,  $I_t^H$  are determined by (also using



that  $p_t^D = p^D$ ):

$$\begin{aligned}
I_t^g + I_t^D + I_t^H + C_t &= L_t w_t + \pi_t + r_t^g K_t^g + r_t^D p_t^D K_t^D + r_t^H p_t^H H_t \\
&\quad + (\tau_t^D (D_t^E + D_t^g) - \tau_t^H p_t^H H_t) \\
I_t^D &\geq 0 \\
I_t^H &\geq 0 \\
I_t^g &= K_{t+1}^g - (1 - \delta^g) K_t^g \\
I_t^D &= p^D (K_{t+1}^D - (1 - \delta^D) K_t^D) \\
I_t^H &= p_t^H (K_{t+1}^H - (1 - \delta^H) K_t^H) \\
\frac{u'(C_t/L_t)}{\beta u'(C_{t+1}/L_{t+1})} &= 1 - \delta^g + r_{t+1}^g \\
r_{t+1}^D &= (1 - \delta^g + r_{t+1}^g) \left( 1 - \frac{\mu_t^{iD}}{u'(C_t/L_t)} \right) - (1 - \delta^D) \left( 1 - \frac{\mu_{t+1}^{iD}}{u'(C_{t+1}/L_{t+1})} \right) \\
\mu_t^{iD} &\geq 0 \\
I_t^D \mu_t^{iD} &= 0 \\
r_{t+1}^H &= \frac{p_t^H}{p_{t+1}^H} (1 - \delta^g + r_{t+1}^g) \left( 1 - \frac{\mu_t^{iH}}{u'(C_t/L_t)} \right) - (1 - \delta^H) \left( 1 - \frac{\mu_{t+1}^{iH}}{u'(C_{t+1}/L_{t+1})} \right) \\
\mu_t^{iH} &\geq 0 \\
I_t^H \mu_t^{iH} &= 0
\end{aligned}$$

## D.1 Compound interest for the firms' problems

Recall our term  $\Pi_t = \Pi_t^g + \Pi_t^D + \Pi_t^H + \Pi_t^{DE} + \Pi_t^E$ . We treated that as a lump-sum. However, in fact the firms are owned by the households, so they choose their activity to maximize the utility pay-off to the households. Thus, for example, the final-goods firms seek to maximize

$$\sum_{t=0}^{\infty} \beta^t \Lambda_t \Pi_t^g$$

subject to its production constraints, where  $\Lambda_t$  is exactly the shadow price on the household's budget constraint above. It is equivalent to divide by  $\Lambda_0$  and so to use a compound discount factor of  $q_t := \beta^t \frac{\Lambda_t}{\Lambda_0} = \beta^t \frac{u'(c_t)}{u'(c_0)}$  for the relative price of consumption in period  $t$ , expressed in period 0 units.

Moreover, recall from (A.74) that  $\frac{\Lambda_t}{\Lambda_{t+1}} = \beta(1 - \delta^g + r_{t+1}^g)$ . Thus

$$\begin{aligned}
q_t &= \beta^t \frac{\Lambda_t}{\Lambda_0} = \frac{\beta \Lambda_t}{\Lambda_{t-1}} \cdot \frac{\beta \Lambda_{t-1}}{\Lambda_{t-2}} \cdots \frac{\beta \Lambda_1}{\Lambda_0} = \prod_{j=1}^t \frac{1}{1 - \delta^g + r_j^g} \\
\frac{q_{t+1}}{q_t} &= \frac{1}{1 - \delta^g + r_{t+1}^g}
\end{aligned} \tag{A.75}$$

## D.2 The final-goods firms' problem

The final-goods firms maximize

$$\sum_{t=0}^{\infty} q_t \Pi_t^g = \sum_{t=0}^{\infty} q_t \left( \Omega(T_t) f(Y_t^g, E_t) - r_t^g K_t^g - w_t L_t - p_t^e E_t - \frac{\phi_{1,t} \eta_t^{\phi_2}}{(1-\eta_t)^{\phi_3}} Y_t^g - p_t^{fuel} D_t^g \right)$$

(remember that  $Y_t^g \equiv f_t^g(K_t^g, L_t)$ ) where  $D_t^g$  are fossil fuels used by these firms,  $p_t^e$  is the price of electricity and  $\frac{\phi_{1,t} \eta_t^{\phi_2}}{(1-\eta_t)^{\phi_3}} Y_t^g$  is spending on abatement by these firms, so that firms face an emission constraint given in every period by:

$$D_t^g = \sigma_t (1 - \eta_t) Y_t^g$$

The first order conditions are then:

$$\partial K_t^g : \quad \Omega(T_t) \frac{\partial f}{\partial Y_t^g} \frac{\partial f_t^g}{\partial K_t^g} = r_t^g + \frac{\phi_{1,t} \eta_t^{\phi_2}}{(1-\eta_t)^{\phi_3}} \frac{\partial f_t^g}{\partial K_t^g} + p_t^{fuel} \sigma_t (1 - \eta_t) \frac{\partial f_t^g}{\partial K_t^g} \quad (\text{A.76})$$

$$\partial L_t : \quad \Omega(T_t) \frac{\partial f}{\partial Y_t^g} \frac{\partial f_t^g}{\partial L_t} = w_t + \frac{\phi_{1,t} \eta_t^{\phi_2}}{(1-\eta_t)^{\phi_3}} \frac{\partial f_t^g}{\partial L_t} + p_t^{fuel} \sigma_t (1 - \eta_t) \frac{\partial f_t^g}{\partial L_t} \quad (\text{A.77})$$

$$\partial E_t : \quad \Omega(T_t) \frac{\partial f_t}{\partial E_t} = p_t^e \quad (\text{A.78})$$

$$\partial \eta_t : \quad p_t^{fuel} \sigma_t = \frac{\phi_{1,t} \eta_t^{\phi_2-1}}{(1-\eta_t)^{1+\phi_3}} [\phi_2 (1 - \eta_t) + \eta_t \phi_3] \quad (\text{A.79})$$

Equation (A.76) is an optimal condition for demand of aggregate capital and states that the return on capital is the marginal product of capital minus additional spending on abatement to clean a given fraction of extra emissions and costs of fuel. Equation (A.76) is the counterpart of equation (A.77) for labor demand. Equation (A.78) is an optimal condition for demand of electricity. Finally, equation (A.79) says that the firm reacts to the price of fuel (implicitly to carbon tax) by choosing the level of abatement (equivalently the level of emissions) such that the price of fuel would be equal to the marginal cost of emissions reduction.

## D.3 Aggregate-electricity-producing firms' problem

The firms produce aggregate electricity by combining both electricity generated by fossil-fuel-based power plants and electricity generated by renewable energy based power stations. Note that we are taking the output from these two plants, in GW, as inputs priced by  $p_t^{EH}$  and  $p_t^{ED}$  respectively and so we do not need to convert by  $p_t^H$  and  $p^D$  here.

$$\sum_{t=0}^{\infty} q_t \Pi_t^E = \sum_{t=0}^{\infty} q_t (p_t^e f_t^E(H_t, \zeta_t K_t^D) - p_t^{EH} H_t - p_t^{ED} (\zeta_t K_t^D))$$

FOCs are:

$$\begin{aligned} p_t^e \frac{\partial f_t^E}{\partial H_t} &= p_t^{EH} \\ p_t^e \frac{\partial f_t^E}{\partial (\zeta_t K_t^D)} &= p_t^{ED} \end{aligned}$$

#### D.4 The dirty-electricity-producing firms' problem

The dirty electricity producing firms are fossil-fuel based power stations, which combine existing infrastructure (for example, coal-based power plants) with fossil fuel, and so maximizes:

$$\sum_{t=0}^{\infty} q_t \Pi_t^D = \sum_{t=0}^{\infty} q_t \left( p_t^{ED} (\zeta_t K_t^D) - r_t^D p^D K_t^D - p_t^{fuel} D_t^E \right)$$

where firms face the emission constraint:  $D_t^E = \nu \zeta_t K_t^D$ , and constraint  $\zeta_t \leq 1$ . So the Lagrangian is (making the obvious substitution)

$$\sum_{t=0}^{\infty} q_t \left( p_t^{ED} (\zeta_t K_t^D) - r_t^D p^D K_t^D - p_t^{fuel} \nu \zeta_t K_t^D + \mu_t^{\zeta} (1 - \zeta_t) \right)$$

And the first order conditions and constraints are

$$\begin{aligned} \partial K_t^D : \quad & r_t^D p^D = \left( p_t^{ED} - p_t^{fuel} \nu \right) \zeta_t \\ \partial \zeta_t : \quad & \mu_t^{\zeta} = K_t^D \left( p_t^{ED} - p_t^{fuel} \nu \right) \\ & \mu_t^{\zeta} (1 - \zeta_t) = 0 \\ & \mu_t^{\zeta} \geq 0 \end{aligned}$$

where  $\mu_t^{\zeta}$  is Lagrangian multiplier attached to the above constraint. Thus, if  $\zeta < 1$  then  $p_t^{ED} = p_t^{fuel} \nu$ , and  $r_t^D p^D = 0$  or  $r_t^D = 0$ . Intuitively, when there is underutilization, the market pushes the return on dirty energy capital to zero.

#### D.5 The fossil-fuel-extracting firm's problem

The firm maximizes

$$\sum_{t=0}^{\infty} q_t \Pi_t^{DE} = \sum_{t=0}^{\infty} q_t [p_t^{fuel} - \tau_t^D - G^D(S_t)] (D_t^E + D_t^g)$$

where  $\tau^D$  is tax on production of fossil fuels. The firm faces the constraint:

$$S_{t+1} = S_t - (D_t^E + D_t^g)$$

to which we assign the shadow price  $\tilde{\mu}_t$ . So the Lagrangian is

$$\mathcal{L}_t = \sum_{t=0}^{\infty} q_t \left( [p_t^{fuel} - \tau_t^D - G^D(S_t)] (D_t^E + D_t^g) - \tilde{\mu}_t (S_{t+1} - S_t + (D_t^E + D_t^g)) \right)$$

FOCs are:

$$\begin{aligned} \partial (D_t^E + D_t^g) : \quad & \tilde{\mu}_t = p_t^{fuel} - \tau_t^D - G^D(S_t) \\ \partial S_{t+1} : \quad & q_t \tilde{\mu}_t = q_{t+1} \left( \tilde{\mu}_{t+1} - (D_{t+1}^E + D_{t+1}^g) (G^D)'(S_{t+1}) \right) \end{aligned}$$

Combining the firm's first order conditions yields the standard Hotelling condition, into which we then substitute from (A.75)

$$\begin{aligned} p_t^{fuel} - \tau_t^D - G^D(S_t) &= \frac{q_{t+1}}{q_t} \left( p_{t+1}^{fuel} - \tau_{t+1}^D - G^D(S_{t+1}) - (D_{t+1}^E + D_{t+1}^g) (G^D)'(S_{t+1}) \right) \\ &= \frac{1}{1 - \delta^g + r_{t+1}^g} \left( p_{t+1}^{fuel} - \tau_{t+1}^D - G^D(S_{t+1}) - (D_{t+1}^E + D_{t+1}^g) (G^D)'(S_{t+1}) \right) \end{aligned}$$

which states that the return on extracting an extra unit of fossil fuels, selling and getting a return on it must be equal to the expected capital gain from keeping an extra unit of fossil fuels in the earth, but extracting it tomorrow minus the increase in future extraction costs. As before, we may repeatedly substitute forward to obtain

$$p_t^{fuel} - \tau_t^D - G^D(S_t) = - \sum_{s=1}^{\infty} \Delta_{t,s} (D_{t+s}^E + D_{t+s}^g) (G^D)'(S_{t+s})$$

where  $\Delta_{t,s} := \prod_{s'=1}^s \frac{1}{1 - \delta^g + r_{t+s'}^g}$ .

## D.6 The renewable energy firms' problem

In contrast to other sectors, we assume that the firms in the renewable sector are small in the sense that they take the stock of accumulated knowledge about using the renewable energy  $H_t$  as given. The renewable energy firms receive subsidy of  $\tau_t^H$  on their dollar-valued holdings of renewable energy capital  $H_t$ . The firms take all prices as given, so they maximize:

$$\sum_{t=0}^{\infty} q_t \Pi_t^H = \sum_{t=0}^{\infty} q_t [p_t^{EH} - p_t^H (r_t^H - \tau_t^H)] H_t.$$

The first order condition is just:

$$p_t^{EH} = p_t^H (r_t^H - \tau_t^H)$$

## D.7 The Principal's Problem

In this section we collect all equations we need to solve the decentralized equilibrium model and formulate it as the principal-agent problem:

$$\max_{\tau^D, \tau^H} \sum_{t=0}^{\infty} \beta^t L_t u \left( \frac{C_t}{L_t} \right)$$

subject to:

$$I_t^g + I_t^D + I_t^H + C_t = L_t w_t + \Pi_t + r_t^g K_t^g + r_t^D p^D K_t^D + r_t^H p_t^H H_t + (\tau_t^D (D_t^E + D_t^g) - \tau_t^H p_t^H H_t) \quad (\text{A.80})$$

$$I_t^D \geq 0$$

$$I_t^H \geq 0$$

$$I_t^g = K_{t+1}^g - (1 - \delta^g) K_t^g$$

$$I_t^D = p^D (K_{t+1}^D - (1 - \delta^D) K_t^D)$$

$$I_t^H = p_t^H (K_{t+1}^H - (1 - \delta^H) K_t^H)$$

$$p_t^H = G(H_t)$$

$$D_t^E = \nu \zeta_t K_t^D$$

$$D_t^g = \sigma_t (1 - \eta_t) Y_t^g$$

$$\frac{u'(C_t/L_t)}{\beta u'(C_{t+1}/L_{t+1})} = 1 - \delta^g + r_{t+1}^g$$

$$r_{t+1}^D = (1 - \delta^g + r_{t+1}^g) \left( 1 - \frac{\mu_t^{iD}}{u'(C_t/L_t)} \right) - (1 - \delta^D) \left( 1 - \frac{\mu_{t+1}^{iD}}{u'(C_{t+1}/L_{t+1})} \right)$$

$$\mu_t^{iD} \geq 0$$

$$I_t^D \mu_t^{iD} = 0$$

$$r_{t+1}^H = \frac{p_t^H}{p_{t+1}^H} (1 - \delta^g + r_{t+1}^g) \left( 1 - \frac{\mu_t^{iH}}{u'(C_t/L_t)} \right) - (1 - \delta^H) \left( 1 - \frac{\mu_{t+1}^{iH}}{u'(C_{t+1}/L_{t+1})} \right) \quad (\text{A.81})$$

$$\mu_t^{iH} \geq 0$$

$$I_t^H \mu_t^{iH} = 0$$

$$\begin{aligned} r_t^g &= \left( \Omega(T_t) \frac{\partial f}{\partial Y_t^g} - \frac{\phi_{1,t} \eta_t^{\phi_2}}{(1 - \eta_t)^{\phi_3}} - p_t^{fuel} \sigma_t (1 - \eta_t) \right) \frac{\partial f_t^g}{\partial K_t^g} \\ &= \left( \Omega(T_t) (1 - \theta) \left[ (1 - \theta) (Y_t^g)^{1-1/\kappa} + \theta (E_t)^{1-1/\kappa} \right]^{\frac{1/\kappa}{1-1/\kappa}} (Y_t^g)^{-\frac{1}{\kappa}} - \frac{\phi_{1,t} \eta_t^{\phi_2}}{(1 - \eta_t)^{\phi_3}} \right. \\ &\quad \left. - p_t^{fuel} \sigma_t (1 - \eta_t) \right) A_t^g \alpha (K_t^g)^{\alpha-1} (L_t)^{1-\alpha} \end{aligned}$$

$$\begin{aligned}
w_t &= \left( \Omega(T_t) \frac{\partial f}{\partial Y_t^g} - \frac{\phi_{1,t} \eta_t^{\phi_2}}{(1 - \eta_t)^{\phi_3}} - p_t^{fuel} \sigma_t (1 - \eta_t) \right) \frac{\partial f_t^g}{\partial L_t} \\
&= \left( \Omega(T_t) (1 - \theta) \left[ (1 - \theta) (Y_t^g)^{1-1/\kappa} + \theta (E_t)^{1-1/\kappa} \right]^{\frac{1/\kappa}{1-1/\kappa}} (Y_t^g)^{-\frac{1}{\kappa}} - \frac{\phi_{1,t} \eta_t^{\phi_2}}{(1 - \eta_t)^{\phi_3}} \right. \\
&\quad \left. - p_t^{fuel} \sigma_t (1 - \eta_t) \right) A_t^g (1 - \alpha) (K_t^g)^\alpha (L_t)^{-\alpha} \tag{A.82}
\end{aligned}$$

$$p_t^e = \Omega(T_t) \frac{\partial f_t}{\partial E_t} = \Omega(T_t) \theta \left[ (1 - \theta) (Y_t^g)^{1-1/\kappa} + \theta (E_t)^{1-1/\kappa} \right]^{\frac{1/\kappa}{1-1/\kappa}} E_t^{-1/\kappa} \tag{A.83}$$

$$\begin{aligned}
p_t^{fuel} \sigma_t &= \frac{\phi_{1,t} \eta_t^{\phi_2-1}}{(1 - \eta_t)^{1+\phi_3}} [\phi_2 (1 - \eta_t) + \eta_t \phi_3] \\
p_t^{EH} &= p_t^e \frac{\partial f_t^E}{\partial H_t} = p_t^e A_t^E \omega H_t^{\xi-1} \left( \omega H_t^\xi + (1 - \omega) (\Gamma_t^{ED})^\xi \right)^{\frac{1-\xi}{\xi}} \tag{A.84}
\end{aligned}$$

$$p_t^{ED} = p_t^e \frac{\partial f_t^E}{\partial (\zeta_t K_t^D)} = p_t^e A_t^E (1 - \omega) (\zeta_t K_t^D)^{\xi-1} \left( \omega H_t^\xi + (1 - \omega) (\Gamma_t^{ED})^\xi \right)^{\frac{1-\xi}{\xi}} \tag{A.85}$$

$$\begin{aligned}
p^D r_t^D &= (p_t^{ED} - p_t^{fuel} \nu) \zeta_t \\
\mu_t^\zeta &= K_t^D (p_t^{ED} - p_t^{fuel} \nu) \tag{A.86}
\end{aligned}$$

$$\begin{aligned}
\mu_t^\zeta (1 - \zeta_t) &= 0 \\
\mu_t^\zeta &\geq 0 \\
p_t^{EH} &= p_t^H (r_t^H - \tau_t^H) \tag{A.87}
\end{aligned}$$

$$p_t^{fuel} - \tau_t^D - G^D(S_t) = - \sum_{s=1}^{\infty} \Delta_{t,s} (D_{t+s}^E + D_{t+s}^g) (G^D)'(S_{t+s}) \tag{A.88}$$

$$\begin{aligned}
\Delta_{t,s} &= \prod_{s'=1}^s \frac{1}{1 - \delta^g + r_{t+s'}^g} \\
D_t &= D_t^E + D_t^{\text{land}} + D_t^g \\
T_t &= \mathcal{W}_t(D_0, \dots, D_{t-1}) \\
S_{t+1} &= S_t - (D_t^E + D_t^g)
\end{aligned}$$

## D.8 Social planner problem versus decentralized equilibrium

**Proof of proposition 5.2** First, from (A.83) and (A.85), we note that:

$$p_t^{ED} = \Omega(T_t) \frac{\partial f_t}{\partial E_t} \frac{\partial f_t^E}{\partial (\zeta_t K_t^D)} = \nu \frac{\partial Y_t}{\partial D_t^E}$$

From (A.86) it follows that:

$$\frac{p^D r_t^D}{\zeta_t \nu} = \frac{p_t^{ED}}{\nu} - p_t^{fuel}$$

And substituting here from the above implies:

$$\frac{\partial Y_t}{\partial D_t^E} = \frac{p^D r_t^D}{\zeta_t \nu} + p_t^{fuel}$$

And substituting the expression for  $p_t^{fuel}$  from (A.88), we obtain:

$$\frac{\partial Y_t}{\partial D_t^E} = \frac{p^D r_t^D}{\zeta_t \nu} + \tau_t^D + G^D(S_t) - \sum_{s=1}^{\infty} \Delta_{t,s} (D_{t+s}^E + D_{t+s}^g) (G^D)'(S_{t+s}) \quad (\text{A.89})$$

Recall that in the social planner's solution, the returns on dirty fuel are equal to (see Proposition C.6):

$$\frac{\partial Y_t}{\partial D_t^E} = \frac{\mu_t^S}{u'(C_t/L_t)} + \chi_t + G^D(S_t) + \frac{p^D R_t^D}{\zeta_t \nu} \quad (\text{A.90})$$

where (see Proposition C.5)

$$\frac{\mu_t^S}{u'(C_t/L_t)} = - \sum_{s=1}^{\infty} \Delta_{t,s} (G^D)'(S_{t+s}) (D_{t+s}^E + D_{t+s}^g)$$

Expression (A.89) is identical to (A.90) when taxes are equal to the SCC, and when  $r_t^D = R_t^D$ .

Next, we find the value of subsidies under which the solutions of the social planner's problem and decentralized equilibrium coincide. First, if the investment into the renewable sector continues then  $\mu_t^{iH} = \mu_{t+1}^{iH} = 0$ , and from (A.81) it follows that:

$$r_{t+1}^H = \frac{p_t^H}{p_{t+1}^H} (1 - \delta^g + r_{t+1}^g) - (1 - \delta^H)$$

or

$$\frac{p_{t+1}^H}{p_t^H} r_{t+1}^H = (1 - \delta^g + r_{t+1}^g) - \frac{p_{t+1}^H}{p_t^H} (1 - \delta^H) \quad (\text{A.91})$$

Using (A.83), (A.84) and (A.87), we can also write that:

$$r_{t+1}^H = \frac{1}{p_{t+1}^H} \frac{\partial Y_{t+1}}{\partial H_{t+1}} + \tau_{t+1}^H \quad (\text{A.92})$$

Next, we denote the return on clean investment in the social planner's case as  $\tilde{r}_{t+1}^H$ . Recall that in the social planner solution (Lemma C.7):

$$\frac{p_{t+1}^H}{p_t^H} \tilde{r}_{t+1}^H = (1 + r_{t+1}^g - \delta^g) - \frac{p_{t+1}^H}{p_t^H} (1 - \delta^H) + \frac{H_{t+2} - (1 - \delta^H) H_{t+1}}{p_t^H} G'(H_{t+1}) \quad (\text{A.93})$$

and

$$\tilde{r}_{t+1}^H = \frac{1}{p_{t+1}^H} \frac{\partial Y_{t+1}}{\partial H_{t+1}} \quad (\text{A.94})$$

Comparison of (A.92) with (A.94) yields the value of subsidies:

$$\tau_{t+1}^H = r_{t+1}^H - \tilde{r}_{t+1}^H$$

But a comparison of (A.91) with (A.93), further yields that:

$$\frac{p_{t+1}^H}{p_t^H} (r_{t+1}^H - \tilde{r}_{t+1}^H) = - \frac{H_{t+2} - (1 - \delta^H) H_{t+1}}{p_t^H} G'(H_{t+1})$$

and the level of subsidies:

$$\tau_t^H = -(H_{t+1} - (1 - \delta^H)H_t) \frac{G'(H_t)}{p_t^H}$$

Finally note that it is straightforward to show that the budget constraint (A.80) is identical to the economy's aggregate constraint as in the social planner's problem after substituting the expressions for profits and returns on capital and labor.  $\square$

## D.9 Discussion of the SCC with the Stringent Damage Factor.

As we mention in the text, our “stringent damage factor”, equation (14), closely approximates a constraint that temperatures cannot exceed the parameter  $\varsigma_5$ . Here we show how the equation for the SCC given in Proposition 5.2 parallels that setting.

By Proposition 5.2, under optimal policy, the carbon tax and SCC are equal to the sum over time of marginal welfare effects from an extra unit of emissions: equation (15). In particular this incorporates the marginal effect of emissions on future output,  $\frac{\partial Y_{t+m}}{\partial D_t}$ , for  $m \geq 1$ . But  $Y_{t+m} = \Omega(T_{t+m})f(Y_{t+m}^g, E_{t+m})$  (equation (7)), and  $T_{t+m} = \mathcal{W}_{t+m}(D_0, \dots, D_{t+m-1})$  (equation (12)). So:

$$\frac{\partial Y_{t+m}}{\partial D_t} = \Omega'(T_{t+m}) \frac{\partial \mathcal{W}_{t+m}}{\partial D_t} f(Y_{t+m}^g, E_{t+m}).$$

Here we unpack the marginal damage factor,  $\Omega'(T_{t+m})$ , in the stringent case.

Write  $\Omega(T_{t+m})$  for the mild damage factor, given in (13), and  $\hat{\Omega}(T_{t+m})$  for the stringent version, given in (14). Recall that  $\hat{\Omega}(T_{t+m}) = \frac{\Omega(T_{t+m})}{1 + \varsigma_3 (T_{t+m}/\varsigma_5)^{\varsigma_4}}$ , where we set  $\varsigma_3$  to be very small ( $\varsigma_3 = 0.001$ ) and  $\varsigma_4$  to be large ( $\varsigma_4 = 50$ ), while  $\varsigma_5$  is the threshold,  $2^\circ\text{C}$ . So the denominator  $1 + \varsigma_3 (T_{t+m}/\varsigma_5)^{\varsigma_4} \approx 1$  when  $T_{t+m} \leq \varsigma_5$ . But this denominator grows rapidly for higher values of  $T_{t+m}$ . It follows that there exists  $\epsilon > 0$  such that  $\hat{\Omega}(T_{t+m}) \approx \Omega(T_{t+m})$  for  $T_{t+m} \leq \varsigma_5$ , and  $\hat{\Omega}(T_{t+m}) \approx 0$  for  $T_{t+m} > \varsigma_5 + \epsilon$ .

It is natural, then, that the mild and stringent marginal damage factors will be approximately equal until  $T_{t+m}$  is just below  $\varsigma_5$ . Meanwhile the marginal stringent damage factor is approximately zero for  $T_{t+m} > \varsigma_5 + \epsilon$ . But  $\hat{\Omega}'(T_{t+m})$  will have to be large and negative at some point between  $\varsigma_5$  and  $\varsigma_5 + \epsilon$ , to allow  $\hat{\Omega}(T_{t+m})$  to move from approximating  $\Omega(T_{t+m})$  (which is close to 1) to being approximately zero.

We indeed observe this when we differentiate (13) and (14) with respect to  $T_{t+m}$ :

$$\begin{aligned} \Omega'(T_{t+m}) &= -\frac{\varsigma_1 \varsigma_2 T_{t+m}^{\varsigma_2-1}}{(1 + \varsigma_1 T_{t+m}^{\varsigma_2})^2} \\ \hat{\Omega}'(T_{t+m}) &= -\frac{\varsigma_1 \varsigma_2 T_{t+m}^{\varsigma_2-1}}{(1 + \varsigma_1 T_{t+m}^{\varsigma_2})^2 (1 + \varsigma_3 (T_{t+m}/\varsigma_5)^{\varsigma_4})^2} - \frac{\varsigma_3 \varsigma_4 (T_{t+m}/\varsigma_5)^{\varsigma_4-1}}{\varsigma_5 (1 + \varsigma_1 T_{t+m}^{\varsigma_2}) (1 + \varsigma_3 (T_{t+m}/\varsigma_5)^{\varsigma_4})^2} \\ &\quad \text{“mild effect”} \qquad \qquad \qquad \text{“threshold effect”} \end{aligned}$$

The marginal stringent damage factor consists of two terms. The first we call the “mild effect”. If  $T_{t+m} \leq \varsigma_5$  then the mild effect  $\approx \Omega'(T_{t+m})$ ; hence the name. Moreover, there exists  $\epsilon_1 > 0$  such that  $(T_{t+m}/\varsigma_5)^{\varsigma_4-1} \approx 0$  for  $T_{t+m} < \varsigma_5 - \epsilon_1$ . Thus, the threshold effect  $\approx 0$  in this range, so that the marginal stringent damage factor is approximated by the mild effect.

There also exists  $\epsilon_2 > 0$  such that  $(T_{t+m}/\varsigma_5)^{\varsigma_4}$  is very large for all  $T_{t+m} > \varsigma_5 + \epsilon_2$ , implying that both the mild and threshold effects are  $\approx 0$  beyond this temperature.

However, if  $T_{t+m}$  is sufficiently close to  $\varsigma_5$ , then the  $\varsigma_4$  (which is large) in the numerator of the threshold effect dominates. So the threshold effect becomes large and negative. We illustrate this



for the parameter values of Table A.1 in Figure A.1. Note that marginal damages below 2 degrees are non-zero, but appear close to zero in comparison with their values just above this temperature. Thus, the contribution to the SCC from each period  $t + m$  breaks down as something close to the DICE damage factor, plus an additional term which becomes very large when we pass the threshold temperature  $\varsigma_5$ .

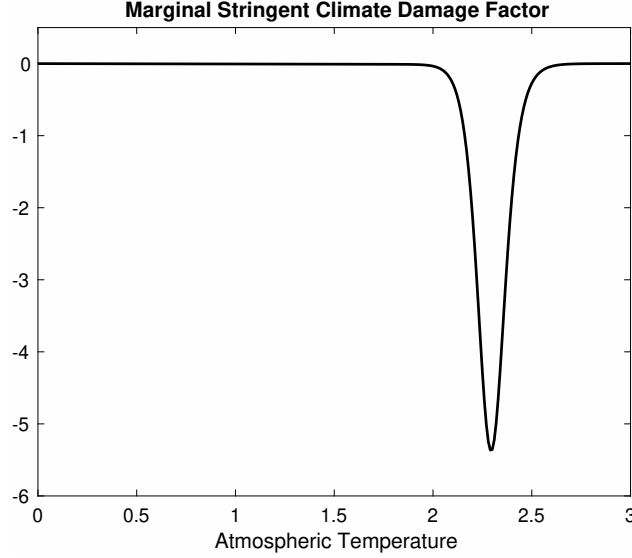


Figure A.1: The Function  $\Omega'(T)$  with  $\varsigma_5 = 2$ .

The precise depth and width of this “downward spike” are determined by our parameters  $\varsigma_3$ ,  $\varsigma_4$  and  $\varsigma_5$ , for which we do not claim empirical justification. The important feature is that the threshold effect, and hence  $\hat{\Omega}'(T_{t+m})$ , depends very sensitively on  $T_{t+m}$  for small changes in  $T_{t+m}$  just above  $\varsigma_5$ . Thus, marginal welfare damages for any given temperature change, range from their value under the mild damage factor, up to very high values, as temperatures increase from just below  $\varsigma_5$  to just above.

The realized marginal welfare losses due to the threshold effect are constrained by the realized temperature. In an optimal trajectory, they will be balanced against marginal economic gains from emitting carbon.

So a small relaxation in the value of  $\varsigma_5$  would have no effect if optimal temperatures never reach this level, and would lead to an increase in welfare commensurate with the avoided cost of keeping to the tighter constraint, if the effective constraint is in fact binding. That is, the threshold effect’s contribution to the SCC is very similar to what we would observe from a shadow price, corresponding to a binding constraint  $T_{t+m} \leq \varsigma_5$ . So, our stringent damage factor can be interpreted as an approximation of this scenario. That is, we can interpret the expression for the SCC as providing, in each period, the sum of marginal damages from our “mild” damage function, plus the shadow cost of keeping temperatures below 2°C in every period.

To further this comparison, consider again the social planner’s problem (Section C), working with the mild damage factor, but impose an additional constraint:

$$T_t \leq \varsigma_5. \quad (\text{A.95})$$

Write  $\mu_t^{max} \geq 0$  for the corresponding period- $t$  shadow price, which is complementary slack with (A.95). Then the Lagrangian (A.41) gains the extra term  $\sum_{t=0}^{\infty} \beta^t \mu_t^{max} (T^{\max} - T_t)$ , so that line

(A.47) is replaced by

$$\partial T_t : \quad \mu_t^W = -\mu_t^{BC} \Omega'(T_t) f_t + \mu_t^{max}. \quad (\text{A.96})$$

Thus, when we derive the SCC in the proof of Proposition C.4 we now substitute (A.96) into (A.46), and divide through by  $\mu_t^{BC}$ , to obtain:

$$\frac{\mu_t^D}{\mu_t^{BC}} = - \sum_{m=1}^{\infty} \beta^m \left( \frac{\mu_{t+m}^{BC}}{\mu_t^{BC}} \Omega'(T_{t+m}) f_{t+m} + \frac{\mu_t^{max}}{\mu_t^{BC}} \right) \frac{\partial \mathcal{W}_{t+m}}{\partial D_t} \quad (\text{A.97})$$

similarly to the previous version. Again, substitute in (A.61) and (A.42), into (A.97) to obtain now:

$$\chi_t := \frac{\mu_t^D}{u'(C_t/L_t)} = \frac{\mu_t^D}{\mu_t^{BC}} = -u' \left( \frac{C_t}{L_t} \right)^{-1} \sum_{m=1}^{\infty} \beta^m u' \left( \frac{C_{t+m}}{L_{t+m}} \right) \left( \frac{\partial Y_{t+m}}{\partial D_t} + \mu_t^{max} \frac{\partial \mathcal{W}_{t+m}}{\partial D_t} \right). \quad (\text{A.98})$$

That is, the social cost of carbon (the shadow price of carbon emissions) is equal to the present value of the marginal economic damages from these emissions in social welfare terms (the first term) plus the marginal effect of these emissions on warming times the shadow price of the temperature constraint (the second term). The latter shadow price is, again, equal to the avoided cost of keeping a slightly tighter constraint. This, then, is exactly what our stringent damage factor approximates.

Moreover, when we solve numerically our social planner's model imposing the constraint (A.95) directly together with the mild damage factor, we find that its solution (including the SCC) is very close to the one using the stringent damage factor without the constraint (A.95). That is, numerically we also show that our stringent damage factor approximates the explicit constraint well.